

*Bernoulli* **25**(3), 2019, 2330–2358  
<https://doi.org/10.3150/18-BEJ1056>

# Bayesian mode and maximum estimation and accelerated rates of contraction

WILLIAM WEIMIN YOO<sup>1</sup> and SUBHASHIS GHOSAL<sup>2</sup>

<sup>1</sup>*Mathematical Institute, Leiden University, P.O. Box 9512, 2300 RA Leiden, The Netherlands.*

*E-mail: [yooweimin0203@gmail.com](mailto:yooweimin0203@gmail.com)*

<sup>2</sup>*Department of Statistics, North Carolina State University, 4276 SAS Hall, 2311 Stinson Drive, Raleigh, NC 27695-8203, USA. E-mail: [sghosal@ncsu.edu](mailto:sghosal@ncsu.edu)*

We study the problem of estimating the mode and maximum of an unknown regression function in the presence of noise. We adopt the Bayesian approach by using tensor-product B-splines and endowing the coefficients with Gaussian priors. In the usual fixed-in-advanced sampling plan, we establish posterior contraction rates for mode and maximum and show that they coincide with the minimax rates for this problem. To quantify estimation uncertainty, we construct credible sets for these two quantities that have high coverage probabilities with optimal sizes. If one is allowed to collect data sequentially, we further propose a Bayesian two-stage estimation procedure, where a second stage posterior is built based on samples collected within a credible set constructed from a first stage posterior. Under appropriate conditions on the radius of this credible set, we can accelerate optimal contraction rates from the fixed-in-advanced setting to the minimax sequential rates. A simulation experiment shows that our Bayesian two-stage procedure outperforms single-stage procedure and also slightly improves upon a non-Bayesian two-stage procedure.

**Keywords:** anisotropic Hölder space; credible set; maximum value; mode; nonparametric regression; posterior contraction; sequential; tensor-product B-splines; two-stage

## 1. Introduction

Consider noisy measurements  $Y_1, \dots, Y_n$  of an unknown smooth function  $f$  at locations  $X_1, \dots, X_n \in [0, 1]^d$  given by the nonparametric regression model

$$Y_i = f(X_i) + \varepsilon_i, \quad i = 1, \dots, n, \quad (1.1)$$

where the regression errors  $\varepsilon_1, \dots, \varepsilon_n$  are modeled as independent and identically distributed (i.i.d.)  $N(0, \sigma^2)$  with unknown standard deviation  $0 < \sigma < \infty$ . The covariates can be deterministic or can be drawn as i.i.d. samples from some fixed distribution independently of the regression errors.

In this paper, we consider the problem of estimating the mode  $\mu$  which marks the location of the maximum of  $f$ , and the value of this maximum  $M = f(\mu) = \sup\{f(x) : x \in [0, 1]^d\}$ , assuming that  $\mu$  is unique. The problem can be thought of as optimization in the presence of noise and has wide range of applications. For instance, searching for the optimal factor configurations in response surface methodology, locating peaks in bacteria (Silverman [26]) and human (Müller [20]) growth curves, or to classify and compare curves arising from longitudinal endocrinological data (Jørgensen et al. [14]).

The problem of estimating the mode and maximum of an isotropic regression function is well studied in the frequentist literature. Müller [20,21] and Shoung and Zhang [25] provided convergence rates for univariate regression, with the multivariate case obtained by Facer and Müller [10]. Furthermore, Hasminskiĭ [13] and Tsybakov [29] showed that for isotropic Hölder regression function of order  $\alpha$  that is also  $\alpha$ -continuously differentiable, the minimax rates for estimating  $\mu$  is  $n^{-(\alpha-1)/(2\alpha+d)}$  and for  $M$  is  $n^{-\alpha/(2\alpha+d)}$ , under the usual sampling plan of choosing samples that are fixed in advance.

However if one is allowed to choose samples based on information gathered from past samples, the structure of the problem changes and we are in the sequential design setting. In this case, the minimax sequential rates of estimating  $\mu$  and  $M$  are respectively  $n^{-(\alpha-1)/(2\alpha)}$  and  $n^{-1/2}$  (see Chen [5], Polyak and Tsybakov [22], Mokkadem and Pelletier [19]). When compared with the fixed design case, it is clear that sequential rates are uniformly better and in fact  $M$  has successfully achieved the parametric rate. Moreover, it also shows that judicious use of past information to guide future actions removes the effect of dimension  $d$  on the rates. On the more practical side, Kiefer and Wolfowitz [15] and Blum [2] used Robbins–Monro type procedures that is consistent; while Fabian [9], Dippon [8] and Mokkadem and Pelletier [19] each constructed sequential procedures that actually attain the minimax rates.

In actual practice, fully sequential design is costly to implement, because sample collection time is longer and the required logistics in collecting data in many stages is much more complicated than single-stage procedures. This then gave rise to the idea of a two-stage procedure, which offers a compromise between the added cost of doing a follow-up experiment and the added accuracy gained from it. At the first stage, limited samples are taken to give a pilot estimate of some quantity (e.g., mode), and the second stage samples are collected in the vicinity of this preliminary estimate. It was then shown in Lan et al. [18], Tang et al. [28] and Belitser et al. [1] that an extra second stage is enough to accelerate the convergence rates and in some cases propel them to attain the minimax sequential rates.

To the best of our knowledge however, there are hardly any such results and procedures in the Bayesian literature, whether it is in the fixed design, sequential or two-stage cases. Therefore, it is hoped that this paper will fill in this gap by giving a Bayesian solution to this problem. As we shall see, there are advantages in using the Bayesian approach, as it provides a natural framework to do two-stage estimation, and it can outperform frequentist procedures by exploiting the shrinkage property of Bayesian estimators.

In the first part of this article, we consider the fixed-in-advance sampling plan and establish single stage posterior contraction rates for  $\mu$  and  $M$ . Our prior consists of tensor product B-splines with Gaussian distributed coefficients, and we endow the error variance with some positive and continuous prior density. We chose this prior because it enables us to derive sharp results by directly analyzing the posterior distribution, and B-splines are efficient to compute. The main challenge here is the non-linear and non-smooth nature of the argmax and max functionals of  $f$ , and we avoid dealing with them directly by relating the estimation errors of  $\mu$  and  $M$  with the sup-norm errors for  $f$  and its first order partial derivatives. To quantify uncertainty in the estimation procedure, we construct credible sets for  $\mu$  and  $M$ , and show that they have high asymptotic coverage with optimal sizes.

Sequential sampling or more specifically a two-stage procedure can naturally be embedded inside a Bayesian framework, as information gained from an earlier stage can be used to adjust

or update one's prior opinion. In the second part of this paper, we propose a Bayesian two-stage procedure for estimating the mode and maximum of  $f$ . We split the samples into two parts, and use the first part to compute the first stage posterior distribution of  $\mu$  and  $M$ . Using this posterior, we construct a credible set based on the techniques discussed in the first part of the paper. Second stage samples are then sampled uniformly over this set, and they are used to compute the second stage posterior of these two quantities.

We show that this second stage posterior is more concentrated around the truth than its single stage counterpart, and it can accelerate single stage minimax rates to the optimal sequential rates, under appropriate conditions on the radius of the credible set used. We test our procedure in a numerical experiment and the results seem to support our theoretical conclusions. Moreover when compared with a non-Bayesian method proposed in the literature, our Bayesian two-stage procedure seems to outperform slightly in terms of the root mean square error, and this is due to the shrinkage induced by our choice of prior distributions (see Figure 3 below).

Throughout this paper, we will work with a general class of anisotropic Hölder space, such that we allow  $f$  to have different order of smoothness in each dimension. In some of our results below, it will be seen that additional smoothness in other dimensions can help alleviate the loss in accuracy due to less smoothness in some dimensions, and this borrowing of smoothness across dimensions, which is a unique feature of anisotropic spaces, can result in the improvement of the overall rate.

The paper is organized as follows. The next section introduces notations and assumptions. Section 3 describes the prior and the resulting posterior distributions of  $\mu$  and  $M$ . Section 4 contains main results in the single stage setting on posterior contraction rates and coverage probability of credible sets for these two quantities. We introduce the Bayesian two-stage procedure of estimating  $\mu$  and  $M$  in Section 5. Section 6 contains simulation studies for our proposed Bayesian two-stage method. This is then followed by a summary and discussion on future outlook in Section 7. Proofs of our main results are given in Section 8 and some useful auxiliary results are collected in the [Appendix](#). We delegate some rather routine and technical proofs to a supplementary article Yoo and Ghosal [34] to streamline reading.

## 2. Notations and assumptions

Given two numerical sequences  $a_n$  and  $b_n$ ,  $a_n = O(b_n)$  or  $a_n \lesssim b_n$  means  $a_n/b_n$  is bounded, while  $a_n = o(b_n)$  or  $a_n \ll b_n$  means  $a_n/b_n \rightarrow 0$ . Also,  $a_n \asymp b_n$  means  $a_n = O(b_n)$  and  $b_n = O(a_n)$ . For stochastic sequence  $Z_n$ ,  $Z_n = O_P(a_n)$  means  $P(|Z_n| \leq Ca_n) \rightarrow 1$  for some constant  $C > 0$ ; while  $Z_n = o_P(a_n)$  means  $Z_n/a_n \rightarrow 0$  in P-probability.

Let  $\|\mathbf{x}\|_p = (\sum_{k=1}^d |x_k|^p)^{1/p}$ ,  $\|\mathbf{x}\|_\infty = \max_{1 \leq k \leq d} |x_k|$  and  $\|\mathbf{x}\| = \|\mathbf{x}\|_2$ . Inequality for a vector stands for co-ordinatewise inequality. For a symmetric matrix  $\mathbf{A}$ , let  $\lambda_{\max}(\mathbf{A})$  and  $\lambda_{\min}(\mathbf{A})$  stand for its largest and smallest eigenvalues, and  $\|\mathbf{A}\|_{(2,2)} = |\lambda_{\max}(\mathbf{A})|$ . Given another matrix  $\mathbf{B}$  of the same size,  $\mathbf{A} \leq \mathbf{B}$  means  $\mathbf{B} - \mathbf{A}$  is nonnegative definite. The  $L_p$ -norm of a function  $f$  is denoted by  $\|f\|_p$ .

We say  $\mathbf{Z} \sim N_J(\boldsymbol{\xi}, \boldsymbol{\Omega})$  if  $\mathbf{Z}$  has a  $J$ -dimensional normal distribution with mean  $\boldsymbol{\xi}$  and covariance matrix  $\boldsymbol{\Omega}$ . By saying that  $\mathbf{Z} \sim \text{GP}(\boldsymbol{\xi}, \boldsymbol{\Omega})$ , we mean that  $\{\mathbf{Z}(t), t \in U\}$  is a Gaussian process with  $E\mathbf{Z}(t) = \boldsymbol{\xi}(t)$  and  $\text{Cov}(\mathbf{Z}(s), \mathbf{Z}(t)) = \boldsymbol{\Omega}(s, t)$  for any  $s, t \in U$ .

Multi-indexes will be frequently used. Let  $\mathbb{N} = \{1, 2, \dots\}$  be the natural numbers and  $\mathbb{N}_0 = \mathbb{N} \cup \{0\}$ . For  $\mathbf{i} = (i_1, \dots, i_d)^T \in \mathbb{N}_0^d$  and  $\mathbf{x} \in \mathbb{R}^d$ , define  $|\mathbf{i}| = \sum_{k=1}^d i_k$ ,  $\mathbf{i}! = \prod_{k=1}^d i_k$  and  $\mathbf{x}^{\mathbf{i}} = \prod_{k=1}^d x_k^{i_k}$ . For  $\mathbf{r} = (r_1, \dots, r_d)^T \in \mathbb{N}_0^d$ , let  $D^{\mathbf{r}} = \partial^{|\mathbf{r}|} / \partial x_1^{r_1} \dots \partial x_d^{r_d}$  be the mixed partial derivative operator. If  $\mathbf{r} = \mathbf{0}$ , we interpret  $D^{\mathbf{0}} f \equiv f$ . If  $\mathbf{r} = \mathbf{e}_k$ , where  $\mathbf{e}_k = (0, \dots, 0, 1, 0, \dots, 0)^T$  with 1 in the  $k$ th position and zero elsewhere, we write  $D^{\mathbf{e}_k}$  as  $D_k$ . We denote  $\nabla f(\mathbf{x}) = (D_1 f(\mathbf{x}), \dots, D_d f(\mathbf{x}))^T$  to be the gradient of  $f$  at  $\mathbf{x}$ . If  $f$  is twice differentiable,  $\mathbf{H} f(\mathbf{x}_0)$  stands for the Hessian matrix of  $f$  at  $\mathbf{x}_0$ .

For  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)^T \in \mathbb{N}^d$ , let us denote  $\alpha^*$  to be the harmonic mean, that is,  $(\alpha^*)^{-1} = d^{-1} \sum_{k=1}^d \alpha_k^{-1}$ . We define the anisotropic Hölder's norm  $\|f\|_{\boldsymbol{\alpha}, \infty}$  as

$$\max \left\{ \|D^{\mathbf{r}} f\|_{\infty} + \sum_{k=1}^d \|D^{(\alpha_k - r_k) \mathbf{e}_k} D^{\mathbf{r}} f\|_{\infty} : \mathbf{r} \in \mathbb{N}_0^d, \sum_{k=1}^d (r_k / \alpha_k) < 1 \right\}. \quad (2.1)$$

The constraint  $\sum_{k=1}^d (r_k / \alpha_k) < 1$  is a technical condition and is imposed so that contraction rates for  $f$  and its derivatives will decrease to 0 as  $n \rightarrow \infty$ .

**Definition 2.1.** The anisotropic Hölder space of order  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_d)^T \in \mathbb{N}^d$ , denoted as  $\mathcal{H}^{\boldsymbol{\alpha}}([0, 1]^d)$ , consists of functions  $f : [0, 1]^d \rightarrow \mathbb{R}$  such that  $\|f\|_{\boldsymbol{\alpha}, \infty} < \infty$ , and for some constant  $C > 0$  with any  $\mathbf{x}, \mathbf{x}_0 \in (0, 1)^d$ ,

$$|D^{\mathbf{r}} f(\mathbf{x}) - D^{\mathbf{r}} T_{\mathbf{x}_0} f(\mathbf{x})| \leq C \sum_{k=1}^d |x_k - x_{0k}|^{\alpha_k - r_k}, \quad (2.2)$$

where  $\mathbf{r} \in \mathbb{N}_0^d$  and  $\sum_{k=1}^d (r_k / \alpha_k) < 1$ . Here  $T_{\mathbf{x}_0} f(\mathbf{x}) = \sum_{\mathbf{i} \leq \mathbf{m}_{\boldsymbol{\alpha}}} D^{\mathbf{i}} f(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0)^{\mathbf{i}} / \mathbf{i}!$  is the tensor Taylor polynomial of order  $\mathbf{m}_{\boldsymbol{\alpha}} := (\alpha_1 - 1, \alpha_2 - 1, \dots, \alpha_d - 1)^T$  by expanding  $f$  around  $\mathbf{x}_0$ .

To study the frequentist properties of the posterior distribution, we assume the existence of a true regression function  $f_0$  such that it satisfies the following three assumptions. In what follows, let  $\mathcal{B}(\mathbf{x}, r) = \{\mathbf{y} : \|\mathbf{y} - \mathbf{x}\| \leq r\}$  be a  $\ell_2$ -ball of radius  $r$  centered at  $\mathbf{x}$ .

1. Under the true distribution  $P_0$ ,  $Y_i = f_0(\mathbf{X}_i) + \varepsilon_i$ , such that  $\varepsilon_i$  are i.i.d. Gaussian with mean 0 and variance  $\sigma_0^2 > 0$  for  $i = 1, \dots, n$ .
2.  $f_0 \in \mathcal{H}^{\boldsymbol{\alpha}}([0, 1]^d)$  for  $\alpha_k > 2, k = 1, \dots, d$ , and attains its maximum  $M_0$  at a unique point  $\boldsymbol{\mu}_0$  in  $(0, 1)^d$  which is well-separated: for any constant  $\tau_1 > 0$ , there exists  $\delta > 0$  such that  $f_0(\boldsymbol{\mu}_0) \geq f_0(\mathbf{x}) + \delta$  for all  $\mathbf{x} \notin \mathcal{B}(\boldsymbol{\mu}_0, \tau_1)$ .
3. For any  $0 < \tau \leq \tau_1$ , there exists  $\lambda_0 > 0$  such that  $\lambda_{\max}\{\mathbf{H} f_0(\mathbf{x})\} < -\lambda_0$  for all  $\mathbf{x} \in \mathcal{B}(\boldsymbol{\mu}_0, \tau)$ .

Assumption 1 states the true regression model for (1.1). The well-separation property of Assumption 2 ensures that only points  $\mathbf{x}$  that are near  $\boldsymbol{\mu}_0$  will give values  $f(\mathbf{x})$  that are close to the true maximum  $M$ . This property is needed to establish posterior consistency for  $\boldsymbol{\mu}$  as we shall see in Theorem 4.1 below. Assumption 3 says that the Hessian of  $f_0$  is locally negative definite around  $\boldsymbol{\mu}_0$ . Observe that Assumptions 2 and 3 imply  $\nabla f_0(\boldsymbol{\mu}_0) = \mathbf{0}$  and the Hessian  $\mathbf{H} f_0(\boldsymbol{\mu}_0)$  is

symmetric and negative definite. Moreover,  $\mathbf{H}f_0(\mathbf{x})$  is continuous in  $\mathbf{x}$ . If  $\alpha_k = 2$ , then we need to make an extra assumption that the second partial derivatives of  $f_0$  are continuous; if not, the Hessian may not be symmetric and its eigenvalues may not be real.

For  $x_k \in [0, 1]$ , let  $B_{j_k, q_k}(x_k)$  be the  $k$ th component B-spline of fixed order  $q_k \geq \alpha_k$ , with knots  $0 = t_{k,0} < t_{k,1} < \dots < t_{k,N_k} < t_{k,N_k+1} = 1$ , such that  $J_k = q_k + N_k$ . Assume that the set of knots in each direction is quasi-uniform, that is,  $\max_{1 \leq l \leq N_k} (t_{k,l} - t_{k,l-1}) \asymp \min_{1 \leq l \leq N_k} (t_{k,l} - t_{k,l-1})$ . Examples include uniform and nested uniform partitions (cf. Examples 6.6 and 6.7 of Schumaker [23]), and we can always choose a subset of knots from any given knot sequence to form a quasi-uniform sequence (cf. Lemma 6.17 of Schumaker [23]).

For fixed design points  $\mathbf{X}_i = (X_{i1}, \dots, X_{id})^T$  with  $i = 1, \dots, n$ , assume that there is a cumulative distribution function  $G$ , with positive and continuous density  $g$  on  $[0, 1]^d$  such that

$$\sup_{\mathbf{x} \in [0, 1]^d} |G_n(\mathbf{x}) - G(\mathbf{x})| = o\left(\prod_{k=1}^d N_k^{-1}\right), \quad (2.3)$$

where  $G_n(\mathbf{x}) = n^{-1} \sum_{i=1}^n \mathbb{1}_{\{\mathbf{X}_i \in [\mathbf{0}, \mathbf{x}]\}}$  is the empirical distribution of  $\{\mathbf{X}_i, i = 1, \dots, n\}$ , with  $\mathbb{1}_U$  the indicator function on  $U$ . The condition holds for the discrete uniform design with  $G$  the uniform distribution when  $N_k \lesssim n^{\alpha^*/\{\alpha_k(2\alpha^*+d)\}}$  for  $k = 1, \dots, d$ . If  $\mathbf{X}_i \stackrel{\text{i.i.d.}}{\sim} G$  with a continuous density on  $[0, 1]^d$ , then (2.3) holds with probability tending to one if  $N_k \lesssim n^{\alpha^*/\{\alpha_k(2\alpha^*+d)\}}$  for  $k = 1, \dots, d$ , and  $\alpha^* > d/2$  by Donsker's theorem. In this paper, we shall prove results on posterior contraction rates and credible sets based on fixed design points. These results will translate to the random case by conditioning on the predictor variables.

### 3. B-splines tensor product, Gaussian prior and posterior

In the model  $Y_i = f(\mathbf{x}_i) + \varepsilon_i$ ,  $i = 1, \dots, n$ , we put a finite random series prior on  $f$  based on tensor-product B-splines, that is,  $f(\mathbf{x}) = \sum_{j_1=1}^{J_1} \dots \sum_{j_d=1}^{J_d} \theta_{j_1, \dots, j_d} \prod_{k=1}^d B_{j_k, q_k}(x_k) := \mathbf{b}_{\mathbf{J}, \mathbf{q}}(\mathbf{x})^T \boldsymbol{\theta}$ , where  $\mathbf{b}_{\mathbf{J}, \mathbf{q}}(\mathbf{x}) = \{\prod_{k=1}^d B_{j_k, q_k}(x_k) : 1 \leq j_k \leq J_k, k = 1, \dots, d\}$  is a collection of  $J = \prod_{k=1}^d J_k$  tensor-product B-splines, and  $\boldsymbol{\theta} = \{\theta_{j_1, \dots, j_d} : 1 \leq j_k \leq J_k, k = 1, \dots, d\}$  are the basis coefficients. Note that  $\mathbf{b}_{\mathbf{J}, \mathbf{q}}(\mathbf{x})$  and  $\boldsymbol{\theta}$  are vectors indexed by  $d$ -dimensional indices and the entries are ordered lexicographically. Then by repeatedly applying equations (15) and (16) of Chapter X from de Boor [7] to each direction  $k = 1, \dots, d$ , the  $\mathbf{r} = (r_1, \dots, r_d)^T$  mixed partial derivative of  $f$  is given by

$$D^{\mathbf{r}} f(\mathbf{x}) = \sum_{j_1=1}^{J_1} \dots \sum_{j_d=1}^{J_d} \theta_{j_1, \dots, j_d} \prod_{k=1}^d \frac{\partial^{r_k}}{\partial x_k^{r_k}} B_{j_k, q_k}(x_k) = \mathbf{b}_{\mathbf{J}, \mathbf{q}-\mathbf{r}}(\mathbf{x})^T \mathbf{W}_{\mathbf{r}} \boldsymbol{\theta}, \quad (3.1)$$

where  $\mathbf{W}_{\mathbf{r}}$  is a  $\prod_{k=1}^d (J_k - r_k) \times \prod_{k=1}^d J_k$  matrix whose entries consist of coefficients associated with applying the finite difference operator iteratively on  $\boldsymbol{\theta}$  (for exact expressions, see (8.1)–(8.4) of Yoo and Ghosal [33]). We represent the model in (1.1) by an  $n$ -variate normal distribution  $\mathbf{Y} | (\mathbf{X}, \boldsymbol{\theta}, \sigma) \sim \mathbf{N}_n(\mathbf{B}\boldsymbol{\theta}, \sigma^2 \mathbf{I}_n)$ , where  $\mathbf{B} = (\mathbf{b}_{\mathbf{J}, \mathbf{q}}(\mathbf{X}_1)^T, \dots, \mathbf{b}_{\mathbf{J}, \mathbf{q}}(\mathbf{X}_n)^T)^T$

is the B-splines basis matrix. Note that we can index the rows and columns of  $\mathbf{B}^T \mathbf{B}$  by multi-dimensional indices, such that for  $\mathbf{u} = (u_1, \dots, u_d)^T$  and  $\mathbf{v} = (v_1, \dots, v_d)^T$ , we write  $(\mathbf{B}^T \mathbf{B})_{\mathbf{u}, \mathbf{v}} = \sum_{i=1}^n \prod_{k=1}^d B_{u_k, q_k}(X_{ik}) B_{v_k, q_k}(X_{ik})$ .

We consider deterministic  $\mathbf{J} = (J_1, \dots, J_d)^T$  number of basis functions depending on  $n, d$  and  $\alpha$ . On the basis coefficients, we endow the prior  $\boldsymbol{\theta} | \sigma^2 \sim \mathcal{N}_J(\boldsymbol{\eta}, \sigma^2 \boldsymbol{\Omega})$ . We choose the prior mean such that  $\|\boldsymbol{\eta}\|_\infty < \infty$  and the  $J \times J$  prior covariance matrix  $c_1 \mathbf{I}_J \leq \boldsymbol{\Omega} \leq c_2 \mathbf{I}_J$  for some constants  $0 < c_1 \leq c_2 < \infty$ . We will use the same multi-dimensional indexing convention of  $\mathbf{B}^T \mathbf{B}$  on  $\boldsymbol{\Omega}$ , and further assume that  $\boldsymbol{\Omega}^{-1}$  is  $\mathbf{h} = (h_1, \dots, h_d)^T$ -banded, in the sense that  $(\boldsymbol{\Omega}^{-1})_{\mathbf{u}, \mathbf{v}} = 0$  if  $|u_k - v_k| > h_k$  for some  $k = 1, \dots, d$ .

By direct calculations,  $D^r f | (\mathbf{Y}, \sigma) \sim \text{GP}(\mathbf{A}_r \mathbf{Y} + \mathbf{c}_r \boldsymbol{\eta}, \sigma^2 \Sigma_r)$ , where  $\mathbf{A}_r$ ,  $\mathbf{c}_r$  and the covariance kernel are defined for  $\mathbf{x}, \mathbf{y} \in (0, 1)^d$  by

$$\mathbf{A}_r(\mathbf{x}) = \mathbf{b}_{J, q-r}(\mathbf{x})^T \mathbf{W}_r (\mathbf{B}^T \mathbf{B} + \boldsymbol{\Omega}^{-1})^{-1} \mathbf{B}^T, \quad (3.2)$$

$$\mathbf{c}_r(\mathbf{x}) = \mathbf{b}_{J, q-r}(\mathbf{x})^T \mathbf{W}_r (\mathbf{B}^T \mathbf{B} + \boldsymbol{\Omega}^{-1})^{-1} \boldsymbol{\Omega}^{-1}, \quad (3.3)$$

$$\Sigma_r(\mathbf{x}, \mathbf{y}) = \mathbf{b}_{J, q-r}(\mathbf{x})^T \mathbf{W}_r (\mathbf{B}^T \mathbf{B} + \boldsymbol{\Omega}^{-1})^{-1} \mathbf{W}_r^T \mathbf{b}_{J, q-r}(\mathbf{y}). \quad (3.4)$$

For  $\sigma$ , we either take an empirical Bayes approach by maximizing the marginal likelihood obtained from  $\mathbf{Y} | \sigma \sim \mathcal{N}_n[\mathbf{B} \boldsymbol{\eta}, \sigma^2 (\mathbf{B} \boldsymbol{\Omega} \mathbf{B}^T + \mathbf{I}_n)]$ , or use a hierarchical Bayes approach. In the former approach, the empirical Bayes estimate given by  $\tilde{\sigma}_n^2 = (\mathbf{Y} - \mathbf{B} \boldsymbol{\eta})^T (\mathbf{B} \boldsymbol{\Omega} \mathbf{B}^T + \mathbf{I}_n)^{-1} \times (\mathbf{Y} - \mathbf{B} \boldsymbol{\eta}) / n$  is plugged into the expression of the conditional posterior process  $D^r f | (\mathbf{Y}, \sigma)$ ; while in the latter approach, we further endow  $\sigma$  with a continuous and positive prior density. For example, we can use a conjugate inverse-gamma (IG) prior  $\sigma^2 \sim \text{IG}(\beta_1/2, \beta_2/2)$ , where  $\beta_1 > 4$  and  $\beta_2 > 0$  are hyper-parameters, to get  $\sigma^2 | \mathbf{Y} \sim \text{IG}[(\beta_1 + n)/2, (\beta_2 + n \tilde{\sigma}_n^2)/2]$ . The posterior for the function  $f$  itself can be recovered as a special case  $\mathbf{r} = \mathbf{0}$  by setting  $\mathbf{W}_0 = \mathbf{I}$ .

**Remark 3.1.** Finite random series based priors have been found to be very convenient to use in Bayesian nonparametrics because their theoretical and computational aspects can be dealt with very simply within the framework of Euclidean spaces. Even though they have simpler structure, they can achieve contraction rates on par with Gaussian process priors. Detailed discussions are given in Shen and Ghosal [24] and the book Ghosal and van der Vaart [12]. More specifically, we used tensor-product B-splines with Gaussian coefficients because it enables us to lower bound the variation of the posterior around its center which is essential to get results on coverage of credible sets (to be discussed in Section 4.2). The structure of the (Gaussian) prior is also helpful for bounding contraction rates and computing the posterior, and this allows us to obtain sharp rates and stronger statements regarding coverage probabilities (e.g. without additional logarithmic factors in the radius), whether it is in the single or two-stage settings. Moreover, B-splines are compactly supported and we have fast recursive algorithm to compute them (see Section 5 of Schumaker [23]). We shall further discuss issues on adaptation in Section 7 where  $\alpha$  is not assumed to be known.

We need to ensure that the priors discussed above for  $f$  will yield a well-defined maximum point at every realization from its posterior. The following lemma assures this property.

**Lemma 3.2.** *If  $f$  is given the tensor-product B-splines with normal coefficients prior, then  $\boldsymbol{\mu}$  is unique for almost all sample paths of  $f$  under the posterior distribution (empirical or hierarchical Bayes).*

In this paper, we simply use  $\Pi(\cdot|Y)$  to denote either the empirical or hierarchical posteriors, we do not distinguish between these two cases since both approaches yield the same rates.

## 4. Main results for single stage setting

### 4.1. Posterior contraction rates

The posterior distributions of  $\boldsymbol{\mu}$  and  $M$  can be induced from the posterior of  $f$  through the argmax and maximum operators. However since these operators are nonlinear and non-differentiable, we take an indirect approach by relating the estimation errors of  $\boldsymbol{\mu}$  and  $M$  to the sup-norm errors in estimating  $f$  and its mixed partial derivatives. By this strategy, results for posterior contraction rates in the supremum norm can be used to induce the corresponding rates for  $\boldsymbol{\mu}$  and  $M$  as the following theorem shows.

**Theorem 4.1.** *Let  $J_k \asymp (n/\log n)^{\alpha^*/\{\alpha_k(2\alpha^*+d)\}}$ ,  $k = 1, \dots, d$  and assume that Assumptions 1, 2 and 3 hold. Consider the empirical Bayes approach by plugging-in  $\tilde{\sigma}_n$  for  $\sigma$ , or the hierarchical Bayes approach by equipping  $\sigma$  with some continuous and positive prior density, then*

$$\|\boldsymbol{\mu} - \boldsymbol{\mu}_0\| \leq \frac{\sqrt{d}}{\lambda_0} \max_{1 \leq k \leq d} \|D_k f - D_k f_0\|_\infty, \quad (4.1)$$

$$|M - M_0| \leq \|f - f_0\|_\infty. \quad (4.2)$$

Therefore for any  $m_n \rightarrow \infty$  and uniformly in  $\|f_0\|_{\alpha, \infty} \leq R$ ,  $R > 0$ ,

$$E_0 \Pi(\|\boldsymbol{\mu} - \boldsymbol{\mu}_0\| > m_n (\log n/n)^{\alpha^* \{1 - (\min_{1 \leq k \leq d} \alpha_k)^{-1}\} / (2\alpha^* + d)} | Y) \rightarrow 0, \quad (4.3)$$

$$E_0 \Pi(|M - M_0| > m_n (\log n/n)^{\alpha^* / (2\alpha^* + d)} | Y) \rightarrow 0. \quad (4.4)$$

Clearly, a consequence of the result above is that the posterior mean  $E(\boldsymbol{\mu}|Y)$  converges to  $\boldsymbol{\mu}_0$  in  $\ell_2$ -norm at the same rate given in (4.3), and the same can be said for  $E(M|Y)$ . Given the absence of minimax results on anisotropic mode estimation (to the best of our knowledge), it is instructive to ask whether the inequalities used above are sharp and the contraction rates obtained are optimal? The following lower bound corollary shows that these results are sharp up to logarithmic factors.

**Corollary 4.2 (Lower bounds).** *In addition to Assumptions 1, 2 and 3, let us now make an extra assumption that for any  $0 < \tau \leq \tau_1$ , we have*

$$\inf_{\mathbf{x}: \|\mathbf{x} - \boldsymbol{\mu}_0\| \leq \tau} \lambda_{\min} \{ \mathbf{H} f_0(\mathbf{x}) \} > -\lambda_1 \quad (4.5)$$



for some constants  $\tau_1, \lambda_1 > 0$ . Then for some small enough constant  $h > 0$ , we have

$$\begin{aligned} E_0 \Pi(\|\boldsymbol{\mu} - \boldsymbol{\mu}_0\| \geq hn^{-\alpha^* \{1 - (\min_{1 \leq k \leq d} \alpha_k)^{-1}\} / (2\alpha^* + d)} \mid \mathbf{Y}) &\rightarrow 1, \\ E_0 \Pi(|M - M_0| \geq hn^{-\alpha^* / (2\alpha^* + d)} \mid \mathbf{Y}) &\rightarrow 1. \end{aligned}$$

For the isotropic smooth case  $\alpha_1 = \dots = \alpha_d = \alpha$ , the norm in (2.1) can be generalized (see Section 2.7.1 of van der Vaart and Wellner [32]) and the B-splines approximation error rate is obtained for all smoothness levels (Theorem 22 of Chapter XII in de Boor [7]). This allows generalization of these and subsequent results in this paper for arbitrary smoothness levels. The contraction rates thus obtained, when reduced to the isotropic case, are the minimax rates for this problem up to some logarithmic factor (see Hasminskii [13] and Tsybakov [29] as discussed in the Introduction). In Section 5 below, we will describe another Bayesian procedure that is able to remove this logarithmic factor and accelerate these rates to the optimal sequential rates.

**Remark 4.3.** In the rates above, clearly the direction which has the least smoothness is the most influential factor in determining the contraction rate for  $\boldsymbol{\mu}$  because of the presence of the second factor in the numerator of the exponent. This is unlike the contraction rate for  $f$ , which is known to be  $(\log n/n)^{\alpha^* / (2\alpha^* + d)}$  (Theorem 4.4 of Yoo and Ghosal [33]), as it depends only on the harmonic mean  $\alpha^*$  of smoothness. The reason is evident from (4.1) in that the largest of the deviations of the function's derivative across all directions bounds the accuracy of estimating  $\boldsymbol{\mu}$ . Nevertheless, it is easily checked that the rate obtained above is better than that obtained by applying the above result on a function of isotropic smoothness  $\min_{1 \leq k \leq d} \alpha_k$ . In other words, additional smoothness in other directions can help to alleviate the comparative loss of accuracy for dimensions which are less smooth, and this borrowing of smoothness across directions, which is a unique feature to anisotropic spaces, results in the improvement of the overall rate.

## 4.2. Credible regions for mode and maximum

Let us now quantify uncertainty by constructing credible regions for  $\boldsymbol{\mu}$  and  $M$ . In what follows, we require that these sets have credibility at least some given level  $1 - \gamma$ , and they have optimal sizes with asymptotic coverage probability also at least  $1 - \gamma$ .

We first construct credible sets in the form of supremum-norm balls in the space of regression functions, and then we map these regions back using the argmax and max functionals into Euclidean spaces, so that they are credible sets for  $\boldsymbol{\mu}$  and  $M$ . Now, the natural Bayesian approach to this problem is to directly construct these sets from the posterior distributions of  $\boldsymbol{\mu}$  and  $M$ . The main reason for favoring the proposed method is that it allows tighter control over the size of the induced credible regions in view of (4.1) and (4.2). Such a control is essential for frequentist coverage, and enables us to use them later in the Bayesian two-stage procedure in Section 5.

We make a remark before we present the main result of this section. It is well known that in nonparametric models, a credible region may have frequentist coverage asymptotically less than the corresponding credibility level. The asymptotic coverage may even be arbitrarily close to zero, because the order of bias of the center of a credible set may be comparable with its variation around the truth under optimal smoothing; see Cox [6], Freedman [11] and Knapik et



al. [17]. This is in sharp contrast with finite dimensional models where Bayesian and frequentist quantification of uncertainty agree because of the Bernstein-von Mises theorem. Knapik et al. [17] showed that this low coverage problem can be addressed by undersmoothing. Castillo and Nickl [3,4] circumvented this problem by using weaker norms to construct credible sets. Szabó et al. [27] and Yoo and Ghosal [33] addressed this same problem by appropriately inflating the size of credible regions to ensure coverage. In our own construction, we shall use the latter approach by introducing a constant  $\rho > 0$  in the radius and choose it large enough so that we will have asymptotic coverage.

Below by posterior distribution we refer to either the empirical Bayes posterior distribution by substituting  $\tilde{\sigma}_n$  for  $\sigma$ , or the hierarchical Bayes posterior distribution obtained by putting a further prior on  $\sigma$ . Denote the posterior mean of  $f$  by  $\tilde{f}$ , and let  $\tilde{\mu}$  be the mode of  $\tilde{f}$  and  $\tilde{M}$  its maximum value. For  $e_k = (0, \dots, 0, 1, 0, \dots, 0)$  with 1 at entry  $k$  and the rest zero, we abbreviate  $A_{e_k}$  by  $A_k$ ,  $c_{e_k}$  by  $c_k$  and  $\Sigma_{e_k}$  as  $\Sigma_k$  respectively, for any  $k = 1, \dots, d$ , where these quantities were defined in (3.2)–(3.4).

**Remark 4.4.** Note that  $\tilde{\mu}$  is well-defined under  $P_0$ . Indeed since the posterior mean is an affine transformation of  $Y$ , it follows from Assumption 1 that  $\tilde{f}$  is a Gaussian process under  $P_0$ . Therefore using the same argument as in the proof of Lemma 3.2, we see that  $\tilde{f}$  has unique maximum  $\tilde{\mu}$  for every realization.

For some  $0 < \gamma < 1/2$ , consider  $\{f : \|D_k f - A_k Y - c_k \eta\|_\infty \leq \rho R_{n,k,\gamma}\}$  as a credible band for  $f$ , where  $\rho > 0$  is a sufficiently large constant and  $R_{n,k,\gamma}$  is the  $(1 - \gamma)$ -quantile of the posterior distribution of  $\|D_k f - A_k Y - c_k \eta\|_\infty$ . Similarly, let  $R_{n,0,\gamma}$  be the  $(1 - \gamma)$ -quantile of the posterior distribution of  $\|f - A_0 Y - c_0 \eta\|_\infty$ . We proceed by using these sets to induce credible regions for  $\mu$  and  $M$  through the argmax and maximum functionals, and they are given by

$$C_\mu = \bigcap_{k=1}^d \{\mu : \|D_k f - A_k Y - c_k \eta\|_\infty \leq \rho R_{n,k,\gamma}\}, \quad (4.6)$$

$$C_M = \{M : \|f - A_0 Y - c_0 \eta\|_\infty \leq \rho R_{n,0,\gamma}\}. \quad (4.7)$$

The following result establishes properties of these regions.

**Theorem 4.5.** *If  $J_k \asymp (n/\log n)^{\alpha^*/\{\alpha_k(2\alpha^*+d)\}}$ ,  $k = 1, \dots, d$ , then we have uniformly in  $\|f_0\|_{\alpha,\infty} \leq R$  for any  $R > 0$ :*

- (i) *the credibility of  $C_\mu$  tends to 1 in  $P_0$ -probability and its coverage approaches 1 asymptotically,*
- (ii)  *$C_\mu \subset \bar{C}_\mu := \{\mu : \|\mu - \tilde{\mu}\|_\infty \leq \rho \sqrt{d} \lambda_0^{-1} \max_{1 \leq k \leq d} R_{n,k,\gamma}\}$  with  $P_0$ -probability going to 1,*
- (iii)  *$\underline{C}_\mu := \{\mu : \|\mu - \tilde{\mu}\|_\infty \leq (Rd)^{-1} \rho \max_{1 \leq k \leq d} R_{n,k,\gamma}\} \subset C_\mu$  with  $P_0$ -probability tending to 1,*
- (iv) *the credibility of  $C_M$  tends to 1 in  $P_0$ -probability and its coverage approaches 1 asymptotically,*

$$(v) \mathcal{C}_M \subset \{M : |M - \tilde{M}| \leq \rho R_{n, \mathbf{0}, \gamma}\}.$$

Assertions (ii) and (iii) say that the induced credible set  $\mathcal{C}_\mu$  can be sandwiched between two hypercubes, and its size is not too small when compared with the upper bound in (ii). Thus, its radius is of the order  $\max_{1 \leq k \leq d} R_{n, k, \gamma} \asymp \max_{1 \leq k \leq d} (\log n/n)^{\alpha^*(1-\alpha_k^{-1})/(2\alpha^*+d)}$  by the second statement of Theorem A.6 in the Appendix (with  $\mathbf{r} = \mathbf{e}_k$ ); while (v) and the same aforementioned statement (with  $\mathbf{r} = \mathbf{0}$ ) imply that the radius of  $\mathcal{C}_M$  is of the order  $(\log n/n)^{\alpha^*/(2\alpha^*+d)}$ . Note that these radius lengths coincide exactly with the contraction rates of Theorem 4.1.

The result above concludes that the induced credible regions for  $\mu$  and  $M$ , that is,  $\mathcal{C}_\mu$  and  $\mathcal{C}_M$  respectively, have adequate frequentist coverage that are of (nearly) optimal sizes. Assertion (ii) also implies that the hypercube  $\bar{\mathcal{C}}_\mu$  centered at  $\tilde{\mu}$  has at least  $(1 - \gamma)$ -credibility and is a confidence set of nearly optimal size. Thus a credible set with guaranteed frequentist coverage can be chosen to be a simple set like a hypercube centered at the posterior mean. In practice, it is easier to construct such a hypercube than the set  $\mathcal{C}_\mu$ , because the latter set requires performing function maximization multiple times to obtain points in  $\mathcal{C}_\mu$ .

The construction of  $\bar{\mathcal{C}}_\mu$  from the data is simple: one finds  $b$  such that the credibility of  $\{\mu : \|\mu - \tilde{\mu}\|_\infty \leq b\}$  is  $1 - \gamma$ , and then inflates this around  $\tilde{\mu}$  by a large constant factor  $\rho > 0$ . For this set to serve as the domain for second stage sampling in a two-stage procedure, some modifications are needed, in that we adjust the length of  $\bar{\mathcal{C}}_\mu$  in each direction so that it adapts to different smoothness. In other words, we embed  $\bar{\mathcal{C}}_\mu$  inside a hyper-rectangle and do uniform sampling inside this larger set. Clearly, keeping the constant inflation factor as small as possible makes the credible sets smaller, but it will be seen that for optimal contraction rate in the second stage, an inflation factor which goes to infinity at a specific rate will be needed.

## 5. Two-stage Bayesian estimation and accelerated rates of contraction

In this section, we show that by obtaining samples in two stages in an appropriate manner, we can accelerate the posterior contraction rates of  $\mu$  and  $M$  to the optimal sequential rates. Given a sampling budget of  $n$ , we first obtain  $n_1 < n$  samples to compute the first stage posterior distribution. The remaining  $n_2 = n - n_1$  samples are then obtained by sampling points uniformly from some regularly shaped credible region constructed from this posterior. Since this is a small region, we can approximate the regression function  $f$  by a multivariate polynomial. By further endowing the coefficients with normal priors, we then use these samples to build the second stage posterior distribution for  $f$  and hence for  $\mu$  and  $M$  through the argmax and maximum functionals. We will then show through Theorem 5.2 below and simulations (Section 6) that these second stage posteriors are more concentrated near the truth.

Let us first describe our proposed Bayesian two-stage procedure in greater detail. Let  $p \in (0, 1)$ . In the first stage, we choose  $n_1 \in \mathbb{N}$  design points  $\{\tilde{\mathbf{x}}_i, i = 1, \dots, n_1\}$  such that  $n_1/n \rightarrow p$  as  $n \rightarrow \infty$  to obtain data  $\mathcal{D}_1 = \{(\tilde{\mathbf{x}}_i, \tilde{Y}_i), i = 1, \dots, n_1\}$  for the model in (1.1). Typically, one chooses  $p = 1/2$  to achieve equal sample splitting but other proportions are possible depending on the sampling configurations (e.g., grid or random sampling) and other practical considerations

such as field conditions and financial constraints. By using the B-spline tensor product prior discussed in Section 3, we obtain a first stage posterior for  $f$ . This then allows us to construct the set

$$\{\boldsymbol{\mu} : |\mu_k - \tilde{\mu}_k| \leq \delta_{n,k}, k = 1, \dots, d\},$$

where  $\tilde{\mu}$  is the mode of the first stage posterior mean of  $f$ . We choose  $\delta_{n,k}, k = 1, \dots, d$  such that

$$\min_{1 \leq k \leq d} \delta_{n,k} = \rho_n \max_{1 \leq k \leq d} (\log n/n)^{\alpha^*(1-\alpha_k^{-1})/(2\alpha^*+d)} \quad (5.1)$$

for a chosen sequence  $\rho_n \rightarrow \infty$ , so that this set is a valid credible set, as it contains  $\bar{\mathcal{C}}_{\boldsymbol{\mu}}$  for large  $n$ . Now sample  $n_2 = n - n_1$  locations  $\{\mathbf{x}_1, \dots, \mathbf{x}_{n_2}\}$  uniformly from this credible set and observe the second stage samples  $\mathcal{D}_2 = \{(\mathbf{x}_i, Y_i), i = 1, \dots, n_2\}$ .

Next, we center the second stage design points at the origin by  $\mathbf{z}_i = \mathbf{x}_i - \tilde{\boldsymbol{\mu}}$  for  $i = 1, \dots, n_2$ . In other words,  $\mathbf{z}_i, i = 1, \dots, n_2$ , are i.i.d. uniform samples from the hyper-rectangle  $\mathcal{Q} := \{\mathbf{x} : |x_k| \leq \delta_{n,k}, k = 1, \dots, d\}$  of sides  $\delta_{n,k}, k = 1, \dots, d$ . Observe that  $\mathbf{z}_i, i = 1, \dots, n_2$ , are independent from the errors  $\boldsymbol{\varepsilon}$  in this sampling scheme. We chose this sampling domain because we need its length at each direction to adapt to different smoothness, and it is operationally more convenient to construct credible sets in the form of hyper-rectangles and do uniform sampling on it (see Remark 5.3 below for a more thorough discussion).

At the second stage, we put a prior on the regression function by representing  $f(\mathbf{z})$  at  $\mathbf{z} = (z_1, \dots, z_d)^T \in \mathcal{Q}$  as a multivariate polynomial function of fixed order  $\mathbf{m}_{\boldsymbol{\alpha}} = (\alpha_1 - 1, \dots, \alpha_d - 1)^T$ , i.e., for  $\mathbf{z}^i = \prod_{k=1}^d z_k^{i_k}$ ,

$$f_{\boldsymbol{\theta}}(\mathbf{z}) = \sum_{i \leq \mathbf{m}_{\boldsymbol{\alpha}}} \theta_i \mathbf{z}^i = \mathbf{p}(\mathbf{z})^T \boldsymbol{\theta}, \quad (5.2)$$

where  $\mathbf{p}(\mathbf{z}) = (\mathbf{z}^i : i \leq \mathbf{m}_{\boldsymbol{\alpha}})^T$  and  $\boldsymbol{\theta} = (\theta_i : i \leq \mathbf{m}_{\boldsymbol{\alpha}})^T$  are the corresponding basis coefficients. The elements of  $\{i : i \leq \mathbf{m}_{\boldsymbol{\alpha}}\}$  can be enumerated as  $\{i_0, i_1, \dots, i_W\}$  where  $W + 1 = \prod_{k=1}^d \alpha_k$  with  $i_0 = \mathbf{0}$ . Define  $\mathbf{Z} = (\mathbf{p}(\mathbf{z}_1), \dots, \mathbf{p}(\mathbf{z}_{n_2}))^T$ , and note that for  $d = 1$ ,  $\mathbf{Z}$  is a Vandermonde matrix.

We endow  $\boldsymbol{\theta}$  with the prior  $\boldsymbol{\theta}|\sigma^2 \sim N_{W+1}(\boldsymbol{\xi}, \sigma^2 \mathbf{V})$ , where the entries of  $\boldsymbol{\xi}$  do not depend on  $n$  and  $\mathbf{V} = \text{diag}\{\prod_{k=1}^d \delta_{n,k}^{-2(i_j)_k} : j = 0, 1, \dots, W\}$ . Then it follows that the posterior  $\Pi(\boldsymbol{\theta}|\mathbf{Y}, \sigma^2)$  is

$$N_{W+1}[(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}(\mathbf{Z}^T \mathbf{Y} + \mathbf{V}^{-1} \boldsymbol{\xi}), \sigma^2(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}]. \quad (5.3)$$

The empirical posterior follows by replacing  $\sigma^2$  with  $\tilde{\sigma}_*^2 = (n_1 \tilde{\sigma}_1^2 + n_2 \tilde{\sigma}_2^2)/n$ , where  $\tilde{\sigma}_1^2 = (\tilde{\mathbf{Y}} - \mathbf{B}\boldsymbol{\eta})^T (\mathbf{B}\boldsymbol{\Omega}\mathbf{B}^T + \mathbf{I}_{n_1})^{-1} (\tilde{\mathbf{Y}} - \mathbf{B}\boldsymbol{\eta})/n_1$  is the empirical estimate of  $\sigma^2$  based on the first stage samples, and  $\tilde{\sigma}_2^2 = n_2^{-1}(\mathbf{Y} - \mathbf{Z}\boldsymbol{\xi})^T (\mathbf{Z}\mathbf{V}\mathbf{Z}^T + \mathbf{I}_{n_2})^{-1} (\mathbf{Y} - \mathbf{Z}\boldsymbol{\xi})$  is the same estimate based on the second stage samples.

For the hierarchical Bayes approach, we use the first stage posterior of  $\sigma^2$  as prior for the second stage. That is, we equip  $\sigma^2$  with  $\text{IG}(\beta_1/2, \beta_2/2)$  prior at the first stage, where  $\beta_1 > 4$

and  $\beta_2 > 0$ , we then use the resulting posterior  $\text{IG}[(\beta_1 + n_1)/2, (\beta_2 + n_1\tilde{\sigma}_1^2)/2]$  as prior for the second stage, which will further yield  $\text{IG}[(\beta_1 + n)/2, (\beta_2 + n\tilde{\sigma}_*^2)/2]$  as the second stage posterior for  $\sigma^2$ . The proposition below shows that the second stage empirical Bayes estimator and the hierarchical Bayes posterior of  $\sigma^2$  are consistent, and it is a key step in establishing the main result given in Theorem 5.2 below.

**Proposition 5.1 (Second stage error variance).** *Uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$ ,*

- (a) *The second stage empirical Bayes estimator  $\tilde{\sigma}_*^2$  converges to  $\sigma_0^2$  in  $P_0$ -probability at the rate  $\max\{n^{-1/2}, n^{-2\alpha^*/(2\alpha^*+d)}, \sum_{k=1}^d \delta_{n,k}^{2\alpha_k}\}$ .*
- (b) *If inverse gamma posterior from the first stage is used as the prior in the second stage, the second stage posterior of  $\sigma^2$  contracts to  $\sigma_0^2$  at the same rate.*

Let  $\mathbf{r} = (r_1, \dots, r_d)^T$  be such that  $\mathbf{r} \leq \mathbf{m}_\alpha$ . Then the  $\mathbf{r}$  mixed partial derivative of  $f_\theta$  is

$$D^{\mathbf{r}} f_\theta(\mathbf{z}) = \sum_{i \leq \mathbf{m}_\alpha} \theta_i \prod_{k=1}^d \frac{\partial^{r_k}}{\partial z_k^{r_k}} z_k^{i_k} = \sum_{\mathbf{r} \leq i \leq \mathbf{m}_\alpha} \theta_i \frac{i!}{(i - \mathbf{r})!} \mathbf{z}^{i - \mathbf{r}}, \quad (5.4)$$

which is a multivariate polynomial of degree  $\mathbf{m}_\alpha - \mathbf{r}$ . The posterior distributions of  $D^{\mathbf{r}} f_\theta$  can then be induced from (5.3).

Let us define the location of the maximum of  $f_\theta$  inside the centered region as  $\mu_z = \arg \max_{\mathbf{z} \in \mathcal{Q}} f_\theta(\mathbf{z})$ . We then relate this location back to the original domain by  $\mu = \tilde{\mu} + \mu_z$  with corresponding maximum value  $M = f_\theta(\mu_z)$ . Following the same reasoning as in Lemma 3.2,  $\mu$  is unique for almost all sample paths of  $f_\theta$  under the empirical or hierarchical posterior. The following theorem establishes the second stage posterior contraction rates of  $\mu$  and  $M$  for any smoothness level  $\alpha_k > 2, k = 1, \dots, d$ , uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$ .

**Theorem 5.2.** *For any chosen sequence  $\rho_n \rightarrow \infty$ , let  $\delta_{n,k}, k = 1, \dots, d$ , be such that  $\min_{1 \leq k \leq d} \delta_{n,k} = \rho_n \max_{1 \leq k \leq d} (\log n/n)^{\alpha_k^*(1-\alpha_k^{-1})/(2\alpha_k^*+d)}$ . Then under Assumptions 1, 2 and 3, we have uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$  and for any  $m_n \rightarrow \infty$ ,*

$$\begin{aligned} \mathbb{E}_0 \Pi \left[ \|\mu - \mu_0\| > m_n \max_{1 \leq k \leq d} \delta_{n,k}^{-1} \left( n^{-1/2} + \sum_{l=1}^d \delta_{n,l}^{\alpha_l} \right) \mid \mathbf{Y} \right] &\rightarrow 0, \\ \mathbb{E}_0 \Pi \left[ |M - M_0| > m_n \left( n^{-1/2} + \sum_{k=1}^d \delta_{n,k}^{\alpha_k} \right) \mid \mathbf{Y} \right] &\rightarrow 0. \end{aligned}$$

*In particular, if  $\alpha_k > 1 + \sqrt{1+d}/2$  for all  $k = 1, \dots, d$ , then for the choice  $\delta_{n,k} = n^{-1/(2\alpha_k)}$ ,  $k = 1, \dots, d$ , the posterior distributions for  $\mu$  and  $M$  contract at the rates  $n^{-(\underline{\alpha}-1)/(2\underline{\alpha})}$  and  $n^{-1/2}$  respectively, where  $\underline{\alpha} = \min_{1 \leq k \leq d} \alpha_k$ .*

Let us take  $\delta_{n,k} = n^{-1/(2\alpha_k)}$  the optimal choice suggested above. By comparing this theorem and the single stage contraction rates of Theorem 4.1, we can draw the following three main conclusions:

1. If we perform a Bayesian two-stage procedure, we accelerate contraction rates for estimating  $\mu$  and  $M$ , with  $M$  achieving the parametric rate  $n^{-1/2}$ .
2. At the same time, we remove extra logarithmic factors that are present in the single stage rates.
3. The second stage rates for  $\mu$  and  $M$  do not depend on  $d$  the dimension of the regression function's domain, and the effect of dimension is mitigated to a lower bound  $1 + \sqrt{1 + d/2}$  required on the smoothness at each direction.

The first conclusion says that the second stage posteriors for  $\mu$  and  $M$  are more concentrated near the truth when compared with their single stage counterparts, and this is evident since the second stage rate of  $\mu$  is  $n^{-(\alpha-1)/(2\alpha)} \ll (\log n/n)^{\alpha^*(1-\alpha^{-1})/(2\alpha^*+d)}$ ; while for  $M$  is  $n^{-1/2} \ll (\log n/n)^{\alpha^*/(2\alpha^*+d)}$ . As noted above,  $M$  has achieved its oracle or the parametric rate. The second stage rate for  $\mu$  is sharp, to see this note that if  $k$  is the worst direction and we know all components of  $\mu$  except the  $k$ th one, then by analyzing the reduced one-dimensional problem, we find the same rate as well since the dimension disappears from the rate. Clearly this is oracle and unbeatable and so the rate is the best possible under anisotropic Hölder spaces. In the isotropic case where  $\alpha_k = \alpha, k = 1, \dots, d$ , the second stage rate for  $\mu$  reduces to  $n^{-(\alpha-1)/(2\alpha)}$  and this is precisely the minimax rate for this problem under sequential sampling (see Chen [5], Polyak and Tsybakov [22] and Belitser et al. [1] as mentioned in the [Introduction](#)).

**Remark 5.3.** There are other ways to construct credible set and obtain the second stage samples. A first attempt would be to simulate  $f$  from its first stage posterior and apply the argmax operator on  $f$ . However as the posterior of  $f$  contracts to  $f_0$ , this approach does not ensure that the second stage samples are sufficiently spread out, that is, the distance between samples at direction  $k$  is at least some constant multiple of  $\delta_{n,k}$  for  $k = 1, \dots, d$ . Another way is to do uniform sampling on  $\mathcal{C}_\mu$  constructed in (4.6), that is, we envelope  $\mathcal{C}_\mu$ , which is possibly irregularly shaped, by the smallest hypercube, do uniform sampling on this cube and discard points that fall outside of  $\mathcal{C}_\mu$ . By (ii) and (iii) of Theorem 5.1,  $\underline{\mathcal{C}}_\mu \subset \mathcal{C}_\mu \subset \overline{\mathcal{C}}_\mu$ , and samples in  $\overline{\mathcal{C}}_\mu$  are proportional to those in  $\underline{\mathcal{C}}_\mu$  under uniform sampling. Hence, the entries of  $(\mathbf{Z}^T \mathbf{Z})^{-1}$  arising from uniform sampling on  $\mathcal{C}_\mu$  or on hypercubes will have the same order, and the second stage posteriors from these two sampling schemes will have the same asymptotic behavior. Operationally, sampling from  $\mathcal{C}_\mu$  requires an extra step in deciding whether the sampled points fall in the set or not.

**Remark 5.4.** For asymptotic analyses, the prior covariance matrix  $\mathbf{V}$  plays a minor role as the data “washes” out the prior. However, for finite samples, the correct specification of  $\mathbf{V}$  is crucial for the success of our proposed method in practical applications. Through empirical experiments, we discovered that  $\mathbf{V}$  must reflect the scaling of the space by  $\delta_{n,k}$  at direction  $k$ , and its inverse must have the same structure as  $\mathbf{Z}^T \mathbf{Z}$  so that  $(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}$  will act as an effective shrinkage factor in (5.3). Let us write  $\mathbf{\Delta} := \text{diag}\{\prod_{k=1}^d \delta_{n,k}^{-(i_j)_k} : j = 0, 1, \dots, W\}$  and we can see that  $\mathbf{Z}^T \mathbf{Z} = \mathbf{\Delta} \mathbf{A} \mathbf{\Delta}$  where  $\mathbf{A}$  is a matrix of constants not depending on  $n$  (see Lemma 8.1 for more details). Therefore if we center the second stage design points, then we match the structure of  $\mathbf{Z}^T \mathbf{Z}$  by choosing  $\mathbf{V} = \mathbf{\Delta}^2$ , which is our default choice. If the points are not centered, we found that Zellner's  $g$ -prior will work, that is,  $\mathbf{V} = g(\mathbf{Z}^T \mathbf{Z})^{-1}$  where  $g$  can be estimated by empirical or hierarchical Bayes methods.

## 6. Simulation study

We shall compare the performance of our two-stage Bayesian procedure with two other estimation methods: the single-stage Bayesian, and the two-stage frequentist procedure proposed in Belitser et al. [1]. Consider the following true regression function defined on  $[0, 1]^2$ :

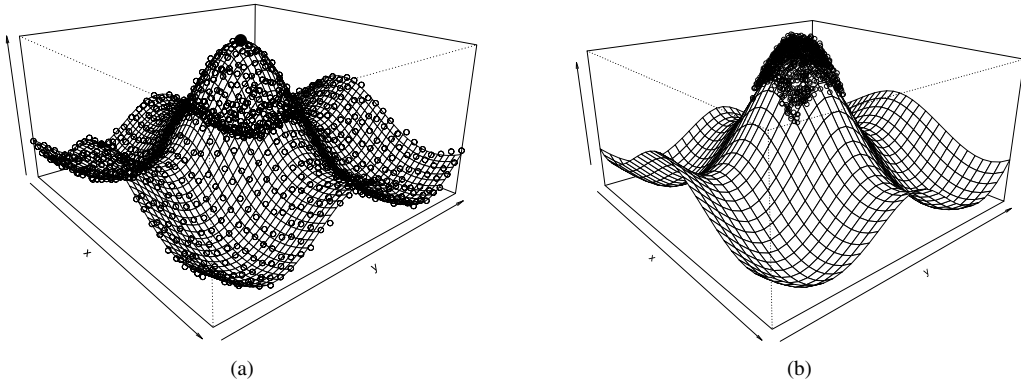
$$f_0(x, y) = (1 + e^{-5(2x-1)^2 - 2(2y-1)^4})[\cos 4(2x - 1) + \cos 5(2y - 1)],$$

where the true mode is given by  $\mu_0 = (0.5, 0.5)^T$ . In the first stage, we observe  $f_0$  on a uniform  $30 \times 30$  grid with i.i.d. errors distributed as  $N(0, 0.01)$  (see Figure 1(a) with black circles as observations). We use bivariate tensor-product of B-splines with normal coefficients as our prior (see Section 3). We choose the pair  $(J_1, J_2)$  that maximizes its posterior, that is,

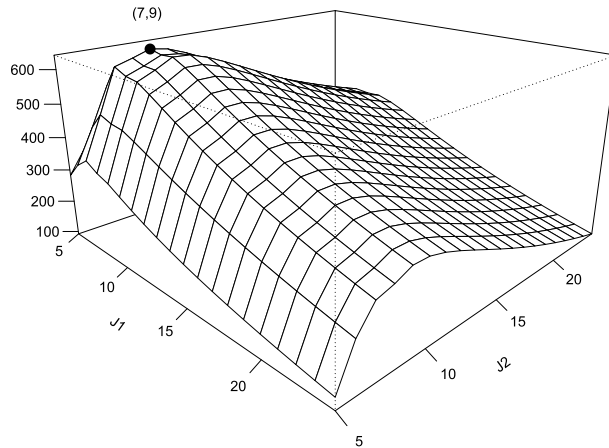
$$\Pi(J_1 = j_1, J_2 = j_2 | Y, \sigma = \tilde{\sigma}_1) \propto \tilde{\sigma}_1^{-n_1} [\det(\mathbf{B}\mathbf{\Omega}\mathbf{B}^T + \mathbf{I}_{n_1})]^{-1/2} \Pi(J_1 = j_1, J_2 = j_2)$$

by integrating out  $\theta$ . We create a candidate set  $\mathcal{J} := \{1, 2, \dots, J_{\max}\} \times \{1, \dots, J_{\max}\}$  by setting  $J_{\max} = 20$ . We put discrete uniform prior on  $(J_1, J_2)$  over  $\mathcal{J}$  such that  $\Pi(J_1 = j_1, J_2 = j_2) = J_{\max}^{-2} = 1/400$  for any  $(j_1, j_2) \in \mathcal{J}$ . We then find the combination that gives the maximum  $\log \Pi(J_1 = j_1, J_2 = j_2 | Y, \sigma = \tilde{\sigma}_1)$  by doing a grid search. To speed up computations, we ignore any constant terms such as the prior factor and the posterior denominator since they do not affect this optimization problem. We plot this marginal log-posterior of  $(J_1, J_2)$  in Figure 2 and we found that  $J_1 = 7$  and  $J_2 = 9$  based on this criterion. At each dimension, the B-spline is of order 4 (cubic) with different uniform knot sequence. For the prior parameters, we set  $\eta = \mathbf{0}$  and  $\mathbf{\Omega} = \mathbf{I}$ . Figure 1(b) shows the surface of the first stage posterior mean of  $f$ .

We sample 864 points uniformly in  $\{\mu : |\mu_k - \tilde{\mu}_k| \leq \delta_k, k = 1, 2\}$  (see black circles in Figure 1(b)) to obtain the second stage data  $Y_i, i = 1, \dots, 864$ . This number is chosen so as to fulfill



**Figure 1.** Our proposed Bayesian two-stage procedure: first stage on the left and second stage on the right. (a) A plot of  $f_0$ , with the black solid point as  $f_0(\mu_0) = M_0$  and black circles as the first stage 900 observations. (b) Posterior mean based on bivariate B-spline prior, and the black circles are the 864 second stage samples.



**Figure 2.** log-posterior of  $(J_1, J_2)$  with its maximum at  $(7, 9)$ .

the sampling configuration required by the frequentist method explained below, and to ensure that all methods under consideration will use the same amount of samples. To choose  $\delta_1, \delta_2$ , we first draw 1000 samples from the first stage posterior distribution of  $f$ , find the mode for each sample by grid search to yield samples from the first stage posterior of  $\mu$ , search for the smallest rectangle enveloping these induced  $\mu$  samples, and take  $\delta_1, \delta_2$  be the lengths of its sides. We found that  $\delta_1 = 0.1111$  and  $\delta_2 = 0.1111$  and we proceed to do uniform sampling across this rectangular region. Note that we do not take the induced  $\mu$  samples as our second stage samples because most of them are concentrated near the center, and hence they are not sufficiently spread out (see Remark 5.3).

We center the 864 sampled design points at the origin by subtracting each of them by  $\tilde{\mu}$ , and we use tensor product of quadratic polynomials with normal coefficients as prior (see (5.2)). That is for  $x, y \in \mathcal{Q} = [-\delta_1, \delta_1] \times [-\delta_2, \delta_2]$ ,  $f_{\theta}(x, y) = \theta_0 + \theta_1 x + \theta_2 y + \theta_3 x^2 + \theta_4 y^2 + \theta_5 xy + \theta_6 xy^2 + \theta_7 x^2 y + \theta_8 x^2 y^2$ . Since  $x, y \in \mathcal{Q}$ , we note that the last three columns of the constructed basis matrix  $\mathbf{Z}$  (corresponding to the last three terms) are very small in magnitude compared with the remaining terms. Thus for numerical simplicity, we consider only

$$f_{\theta}^*(x, y) = \theta_0 + \theta_1 x + \theta_2 y + \theta_3 x^2 + \theta_4 y^2 + \theta_5 xy,$$

where  $\theta | \sigma^2 \sim N(\mathbf{0}, \sigma^2 \mathbf{V})$ , with  $\mathbf{V} = \text{diag}(1, \delta_1^{-2}, \delta_2^{-2}, \delta_1^{-4}, \delta_2^{-4}, \delta_1^{-2} \delta_2^{-2})$ . We use the empirical Bayes method to estimate  $\sigma$  by  $\tilde{\sigma}_*$ , which is the weighted average of the first and second stage estimates  $\tilde{\sigma}_1, \tilde{\sigma}_2$ . However, we note that in our simulations that  $\tilde{\sigma}_2$  gives a much better estimate than  $\tilde{\sigma}_1$  at the current  $(n_1, n_2)$ -sampling plan, and since both are valid independent estimates of  $\sigma$ , we take  $\tilde{\sigma}_* = \tilde{\sigma}_2$ . Now, to compute  $\mu_z = \arg \max_{z \in \mathcal{Q}} f_{\theta}^*(z)$  for a fixed  $\theta$ , we solve the following system of equation  $\nabla f_{\theta}^*(\mu_z) = \mathbf{0}$ , which is equivalent to solving

$$\begin{pmatrix} 2\theta_3 & \theta_5 \\ \theta_5 & 2\theta_4 \end{pmatrix} \begin{pmatrix} \mu_{z,1} \\ \mu_{z,2} \end{pmatrix} = \begin{pmatrix} -\theta_1 \\ -\theta_2 \end{pmatrix} \quad (6.1)$$



for  $\mu_z = (\mu_{z,1}, \mu_{z,2})^T$ . Therefore, to induce the posterior distribution of  $\mu$ , we draw samples from  $\Pi(\theta|Y)$  by substituting  $\sigma = \tilde{\sigma}_2$  in (5.3), solving for  $\mu_z$  using (6.1) for each sample, and translating back to  $\mu = \tilde{\mu} + \mu_z$ .

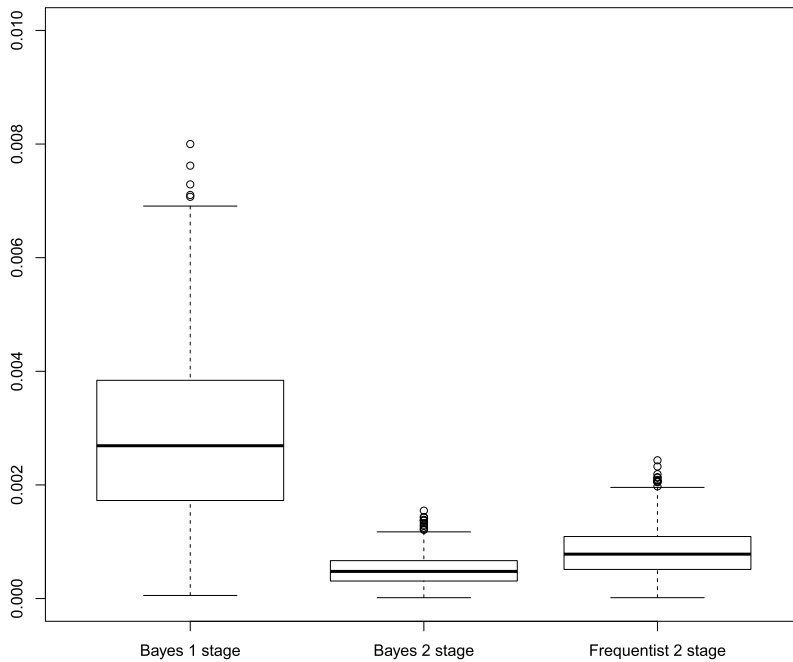
The procedure described is implemented in the statistical software package R. Univariate B-splines are constructed using the `bs` function from the `splines` package. We then use `tensor.prod.model.matrix` from the `crs` package to form their tensor products. Observe that we have used a total of  $30 \times 30 + 864 = 1764$  observations. To show that the two-stage procedure indeed has better accuracy than single stage procedures, we then compare our two-stage procedure with a Bayesian single stage method that uses the same number of samples, that is, on a uniform  $42 \times 42$  grid points. As in the previous case, we observe  $f_0$  at these points with i.i.d. errors  $N(0, 0.01)$ . We then use bivariate tensor-product B-splines with normal coefficients as prior with the same setting. We found that  $J_1 = 9, J_2 = 9$  maximizes its posterior.

To compare Bayesian and frequentist procedures, we repeated the same experiment using the two-stage frequentist procedure implemented in Belitser et al. [1]. In their procedure, we first fit a locally linear surface by `loess` regression in R on the first stage design points. We choose the corresponding bandwidth or span parameter to be 0.02 by leave-one-out cross validation. The maximum of this surface serves as a preliminary estimator. We construct a rectangle of size  $2\delta_1 \times 2\delta_2$  around this estimator and further divide this region into 4 smaller  $\delta_1 \times \delta_2$  rectangles. In their implementation, Belitser et al. [1] actually tuned  $\delta_1, \delta_2$  using the knowledge of the true maximum, and to the best of our knowledge, there is no practical frequentist method to choose the  $\delta$ s in the literature. Faced with this situation, we decided to follow Belitser et al. [1] and choose the  $\delta$ s by minimizing the expected  $L_2$ -distance between the second stage posterior mean for  $\mu$  and the true  $\mu_0$ . We found that  $\delta_1 = 0.06$  and  $\delta_2 = 0.06$ . Then, we take 96 replicated samples at the 9 grid points to form 864 second stage samples. A quadratic surface is fit through these points and its coefficients are estimated using least squares. We compute the second stage estimator of  $\mu_0$  by (6.1) and call it  $\hat{\mu}_2$ .

Let  $\tilde{\mu}_1$  and  $\tilde{\mu}_2$  be the single stage and two-stage posterior means respectively. To compare the performance of our proposed Bayesian two-stage method with the other two, we replicate the experiment 1000 times for each of the three methods. For each replicated experiment, we compute  $\|\mu - \mu_0\|$  for each  $\mu = \tilde{\mu}_1, \tilde{\mu}_2, \hat{\mu}_2$ . Figure 3 shows the box-plots of these 1000 computed root mean square errors (RMSE) for all three procedures.

We see that the our proposed Bayesian two-stage procedure has considerably lower RMSE than the corresponding single stage procedure, and thus supports the conclusion of Theorem 5.2 in finite sample setting. We also observe that our proposed Bayesian two-stage procedure performs slightly better than the frequentist procedure, despite the fact that we have used the true maximum to tune the frequentist procedure while our proposed Bayesian two-stage method is fully data driven. The superior performance may be due to our choice of prior covariance matrix  $V$  (see Remark 5.4) where it causes  $(Z^T Z + V^{-1})^{-1}$  in the posterior (see (5.3)) to act as an effective shrinkage factor.

The R codes to reproduce all the results and figures in this section can be found in the first author's GitHub <https://github.com/wwyoo/Bayes-two-stage>.



**Figure 3.** Root mean square errors (y-axis) for the Bayesian and frequentist procedures with 1000 Monte Carlo replications.

## 7. Conclusion and future works

We studied Bayesian estimation of  $\mu$  and  $M$  in two different settings. In the usual single stage situation, we obtained posterior contraction rates of  $(\log n/n)^{\alpha^*\{1-\alpha^{-1}\}/(2\alpha^*+d)}$  for  $\mu$  and  $(\log n/n)^{\alpha^*/(2\alpha^*+d)}$  for  $M$ , under a tensor-product B-splines random series prior with Gaussian coefficients. Using our proposed Bayesian two-stage procedure, we can accelerate the aforementioned rates to  $n^{-(\alpha-1)/(2\alpha)}$  and  $n^{-1/2}$  for  $\mu$  and  $M$  respectively, as long as the optimal  $\delta_{n,k} = n^{-1/(2\alpha_k)}$  is chosen as radius for the credible cube used in second stage sampling. This rate acceleration is remarkable because it removes the logarithmic factors and it mitigates the effect of dimension  $d$  on the rates. We implemented a practical version of our Bayesian two-stage procedure in a simulation study, and it outperformed a traditional single stage Bayesian procedure by a large margin, and also performed slightly better compared to a frequentist procedure recently proposed in the literature.

An important future work for the single and two-stage Bayesian procedures is to make them adaptive to the unknown smoothness  $\alpha$ . In other words, designing theoretically sound and data driven procedures to determine the optimal  $J_k$  (number of B-splines) and  $\delta_{n,k}$  (credible cube radius) for  $k = 1, \dots, d$ . If only  $L_2$ -distances are studied, finite random series easily gives adaptation by simply putting a prior on  $J_k$  the number of basis functions. For supremum  $L_\infty$ -distance as considered in this paper, getting adaptive posterior contraction rate is a lot more challenging

and seems to need a different type of prior (see Yoo et al. [35]). Coverage of uniform norm credible sets in the adaptive setting may even need a more radical technique (cf. Yoo and van der Vaart [36]).

For two-stage procedures, rate adaptation has not been yet possible even for the non-Bayesian procedure of Belitser et al. [1]. At the minimum, to adapt, one may need multi-stage (possibly more than 2) sampling. In our simulation study, the empirical Bayes step used in determining  $J_k$  and the sampling based procedure to choose  $\delta_1, \delta_2$  are practical and fast plug-in methods, but it is yet unknown how estimation uncertainty introduced by plugging-in these estimated quantities propagate throughout the two-stage procedure, and whether this is an optimal thing to do, even though they do give reasonable finite sample results. We shall try to answer these pressing questions in some future work.

## 8. Proofs

Recall that the posterior mean of  $D^r f$  is  $D^r \tilde{f} := A_r Y + c_r \eta$ . Let us write the posterior contraction rate of  $\mu$  given in Theorem 4.1 as  $\epsilon_n = \max_{1 \leq k \leq d} \epsilon_{n,k}$ , where  $\epsilon_{n,k} = (\log n/n)^{\alpha^*(1-\alpha_k^{-1})/(2\alpha^*+d)}$ .

**Proof of Lemma 3.2.** For the empirical Bayes case, the reproducing kernel Hilbert space of the Gaussian posterior process  $\{f(t) : t \in [0, 1]^d\}$  given  $Y$  is the  $J$ -dimensional space of polynomial splines spanned by elements of  $b_{J,q}(\cdot)$  (see Theorem 4.2 of van der Vaart and van Zanten [31]). This implies that it has continuous sample path at every realization. Observe that  $\tilde{\sigma}_n^2 > 0$  almost surely. Then if

$$\begin{aligned} 0 &= \text{Var}(f(t) - f(s) | Y, \sigma = \tilde{\sigma}_n) \\ &= \tilde{\sigma}_n^2 (b_{J,q}(t) - b_{J,q}(s))^T (B^T B + \Omega^{-1})^{-1} (b_{J,q}(t) - b_{J,q}(s)), \end{aligned}$$

this implies that  $B_{j_k, q_k}(t_k) = B_{j_k, q_k}(s_k)$  for  $j_k = 1, \dots, J_k$ , and  $k = 1, \dots, d$ . Since  $q_k \geq \alpha_k > 2$  for  $k = 1, \dots, d$ , this rules out the possibility that the B-splines are step functions and further implies that  $t = s$ . Thus by Lemma 2.6 of Kim and Pollard [16],  $\mu$  is unique for almost all sample paths of  $\{f(t) : t \in [0, 1]^d\}$  under the empirical posterior distribution.

For hierarchical Bayes, consider the conditional posterior process  $\Pi(f|Y, \sigma)$  for arbitrary  $\sigma > 0$ . The reproducing kernel Hilbert space of this Gaussian process is still the space of splines as before, and consequently has continuous sample path for each realization. Now by substituting  $\sigma$  for  $\tilde{\sigma}_n$  in the preceding display, we see that  $\text{Var}(f(t) - f(s) | Y, \sigma) = 0$  implies  $t = s$ . Again by Lemma 2.6 of Kim and Pollard [16], almost every sample path of  $f$  given  $\sigma$  has unique maximum for any  $\sigma > 0$ . Note that  $f$  is generated from the following scheme: draw  $\sigma \sim \Pi(\sigma|Y)$ , and then  $f|\sigma \sim \text{GP}(AY + c\eta, \sigma^2 \Sigma)$ . Hence almost every draw of  $f$  has unique maximum  $\mu$  under  $\Pi(f|Y)$ .  $\square$

**Proof of Theorem 4.1.** If both  $f, f_0 \geq 0$ , then (4.2) follows from the reverse triangle inequality by noting that  $M = \|f\|_\infty$  and  $M_0 = \|f_0\|_\infty$  in this case. If the condition fails to hold, we can

add a large enough constant  $C > 0$  such that  $g = f + C \geq 0$  and  $g_0 = f_0 + C \geq 0$ . Then by the reverse triangle inequality,  $|\|g\|_\infty - \|g_0\|_\infty| \leq \|g - g_0\|_\infty$ . The right-hand side is  $\|f - f_0\|_\infty$ , while the left-hand side is  $|M - M_0|$  because  $\|g\|_\infty = M + C$ ,  $\|g_0\|_\infty = M_0 + C$ .

To prove (4.1), we need to first establish consistency of the induced posterior of  $\mu$ . Now for any  $\epsilon > 0$  and  $\delta > 0$ ,  $\Pi(\|\mu - \mu_0\| > \epsilon \mid Y)$  is bounded above by

$$\begin{aligned} \Pi\left(\sup_{x \notin \mathcal{B}(\mu_0, \epsilon)} f(x) > f(\mu_0) \mid Y\right) &\leq \Pi\left(\sup_{x \notin \mathcal{B}(\mu_0, \epsilon)} f(x) > f(\mu_0), f(\mu_0) \geq f_0(\mu_0) - \delta/2 \mid Y\right) \\ &\quad + \Pi(f(\mu_0) < f_0(\mu_0) - \delta/2 \mid Y). \end{aligned}$$

The second term is bounded above by  $\Pi(|f(\mu_0) - f_0(\mu_0)| > \delta/2 \mid Y) \leq \Pi(\|f - f_0\|_\infty > \delta/2 \mid Y)$ , and this goes to 0 in  $P_0$ -probability by Theorem A.5 in the Appendix with  $r = \mathbf{0}$ . The well-separation property of Assumption 2 implies that there exists a  $\delta > 0$ , such that  $f_0(x) < f_0(\mu_0) - \delta$  for  $x \notin \mathcal{B}(\mu_0, \epsilon)$ . Hence, for this  $\delta > 0$  and appealing again to Theorem A.5, the first term is bounded above by

$$\Pi\left(\bigcup_{x \notin \mathcal{B}(\mu_0, \epsilon)} \left\{f(x) > f_0(x) + \frac{\delta}{2}\right\} \mid Y\right) \leq \Pi\left(\|f - f_0\|_\infty > \frac{\delta}{2} \mid Y\right) \rightarrow 0$$

in  $P_0$ -probability as  $n \rightarrow \infty$ . Thus the posterior distribution of  $\mu$  is consistent at  $\mu_0$ .

Now by Taylor's theorem,  $\nabla f_0(\mu) = \nabla f_0(\mu_0) + \mathbf{H} f_0(\mu^*)(\mu - \mu_0)$  for some  $\mu^* = \lambda\mu + (1 - \lambda)\mu_0$  with  $\lambda \in (0, 1)$ . Since the posterior of  $\mu$  is consistent as shown above and  $\mu^*$  falls in between  $\mu$  and  $\mu_0$ , it must be that as  $n \rightarrow \infty$  and for any  $\epsilon > 0$ ,  $\Pi(\|\mu^* - \mu_0\| \leq \epsilon \mid Y) \xrightarrow{P_0} 1$ . Let us introduce the set  $\mathcal{B} = \{\lambda_{\max}[\mathbf{H} f_0(\mu^*)] < -\lambda_0\}$ , and note that under Assumptions 3,  $\Pi(\mathcal{B} \mid Y) \xrightarrow{P_0} 1$  and  $\mathbf{H} f_0(\mu^*)$  is invertible with posterior probability tending to one. Therefore, as  $n \rightarrow \infty$  and intersecting with  $\mathcal{B}$ ,  $\mu - \mu_0 = \mathbf{H} f_0(\mu^*)^{-1}(\nabla f_0(\mu) - \nabla f_0(\mu_0))$ . Noting that  $\nabla f_0(\mu_0) = \nabla f(\mu) = \mathbf{0}$  by Assumption 2, then

$$\|\mu - \mu_0\| \leq \frac{1}{\lambda_0} \|\nabla f_0(\mu) - \nabla f(\mu)\| \leq \frac{\sqrt{d}}{\lambda_0} \max_{1 \leq k \leq d} \|D_k f - D_k f_0\|_\infty,$$

where we have used the sub-multiplicative property of the  $\|\cdot\|_{(2,2)}$ -norm, that is,  $\|A\mathbf{y}\| \leq \|A\|_{(2,2)} \|\mathbf{y}\|$  for some matrix  $A$  and vector  $\mathbf{y}$ .

Theorem A.5 with  $r = \mathbf{0}$  together with (4.2) now proves the desired contraction rate on  $M$ . To derive the rate (4.3) for  $\mu$ , apply again Theorem A.5 with  $r = \mathbf{e}_k$ , and the  $L_\infty$ -contraction rate for  $D_k f$  is  $\epsilon_{n,k}$ ,  $k = 1, \dots, d$ . In view of (4.1), we have for any  $m_n \rightarrow \infty$ ,

$$\mathbb{E}_0 \Pi(\|\mu - \mu_0\| > m_n \epsilon_n \mid Y) \leq \sum_{k=1}^d \mathbb{E}_0 \Pi\left(\|D_k f - D_k f_0\|_\infty > \frac{\lambda_0}{\sqrt{d}} m_n \epsilon_n \mid Y\right)$$

approaches 0 uniformly in  $\|f_0\|_{\alpha, \infty} \leq R$ , establishing the assertion. □

**Proof of Corollary 4.2.** See supplementary article Yoo and Ghosal [34]. □

**Proof of Theorem 4.5.** By construction,  $\mathcal{C}_\mu$  contains  $\tilde{\mu}$ , the mode of the posterior mean  $\tilde{f}$  of  $f$ . Since  $\gamma < 1/2$ ,  $R_{n,k,\gamma}$  is greater than the posterior median of  $\|D_k f - D_k \tilde{f}\|_\infty$ . For the empirical Bayesian posterior,  $D_k f - D_k \tilde{f}$  is a Gaussian process under Assumption 1 and in view of the second assertion of Theorem A.6 in the Appendix, it follows that  $R_{n,k,\gamma} \gtrsim \epsilon_{n,k}$  (with  $\mathbf{r} = \mathbf{e}_k$ ). Define  $v_{n,k}^2 := \sup_{\mathbf{x} \in [0,1]^d} \text{var}(D_k f(\mathbf{x}) - D_k \tilde{f}(\mathbf{x}) | \mathbf{Y})$ . Since the empirical Bayes estimate  $\tilde{\sigma}_n^2$  is consistent in view of (a) in Proposition A.9, then by applying the inequality  $\mathbf{y}^T \mathbf{A} \mathbf{y} \leq \lambda_{\max}(\mathbf{A}) \|\mathbf{y}\|^2$  for any square matrix  $\mathbf{A}$ , we can bound  $\text{var}(D_k f(\mathbf{x}) - D_k \tilde{f}(\mathbf{x}) | \mathbf{Y})$  with expression given in (3.4) by

$$\begin{aligned} & (\sigma_0^2 + o_{P_0}(1)) \lambda_{\min}^{-1}(\mathbf{B}^T \mathbf{B} + \boldsymbol{\Omega}^{-1}) \lambda_{\max}(\mathbf{W}_{\mathbf{e}_k} \mathbf{W}_{\mathbf{e}_k}^T) \|\mathbf{b}_{\mathbf{J}, \mathbf{q} - \mathbf{e}_k}(\mathbf{x})\|^2 \\ & \lesssim n^{-1} J_k^2 \prod_{l=1}^d J_l \lesssim n^{-2\alpha^*(1-\alpha_k^{-1})/(2\alpha^*+d)} \end{aligned}$$

for any  $\mathbf{x} \in [0,1]^d$ . In the above, we have used the fact  $0 \leq B_{j_l, q_l}(x_l) \leq 1$  and the partition of unity property of B-splines  $\sum_j \prod_{l=1}^d B_{j_l, q_l}(x_l) = 1$  to bound  $\|\mathbf{b}_{\mathbf{J}, \mathbf{q} - \mathbf{e}_k}(\mathbf{x})\|^2 \leq \sum_{j_1=1}^{J_1} \cdots \sum_{j_d=1}^{J_d} \prod_{l=1}^d B_{j_l, q_l - \mathbb{1}_{\{l=k\}}}(x_l) \leq 1$  for any  $\mathbf{x} \in [0,1]^d$ . The eigenvalues were bounded by (A.13) of Lemma A.7 and Lemma A.8 (with  $\mathbf{r} = \mathbf{e}_k$ ). Thus, we conclude that  $v_{n,k}^2 = o(\epsilon_{n,k}^2)$ . Consequently by Borell's inequality (cf. Proposition A.2.1 of van der Vaart and Wellner [32]) and taking  $\rho$  to be large enough, for example,  $\rho > 1$ ,

$$\begin{aligned} \Pi(\mu \notin \mathcal{C}_\mu | \mathbf{Y}) & \leq \sum_{k=1}^d \Pi(\|D_k f - D_k \tilde{f}\|_\infty > \rho R_{n,k,\gamma} | \mathbf{Y}) \\ & \leq d \max_{1 \leq k \leq d} \exp[-(\rho - 1)^2 R_{n,k,\gamma}^2 / (2v_{n,k}^2)] \end{aligned}$$

which goes to zero since  $R_{n,k,\gamma}^2 / v_{n,k}^2 \rightarrow \infty$ . Thus, the credibility of  $\mathcal{C}_\mu$  tends to 1 (or is at least  $1 - \gamma$ ) in  $P_0$ -probability as  $n \rightarrow \infty$ . For the hierarchical Bayesian posterior, the same conclusion follows since conditionally on  $\sigma$ , the posterior law obeys a Gaussian process and  $\sigma$  lies in a small neighborhood of the true  $\sigma_0$  with high posterior probability (from (c) of Proposition A.9). To ensure coverage, note that by the construction of  $\mathcal{C}_\mu$ , it contains  $\mu_0$  if  $\|D_k f_0 - D_k \tilde{f}\|_\infty \leq \rho R_{n,k,\gamma}$  for all  $k = 1, \dots, d$ . When  $f_0$  is the true regression function, the  $P_0$ -probability of the last event tends to one for all  $k$  by Theorem A.6 with  $\mathbf{r} = \mathbf{e}_k$ , and hence the statement on coverage is established. This proves assertion (i).

Assertion (ii) follows in view of (4.6) and if

$$\|\tilde{\mu} - \mu_0\| \leq \frac{\sqrt{d}}{\lambda_0} \max_{1 \leq k \leq d} \|D_k \tilde{f} - D_k f_0\|_\infty \quad (8.1)$$

holds with  $P_0$ -probability tending to 1. Indeed by Remark 4.4,  $\nabla \tilde{f}(\tilde{\mu}) = \mathbf{0}$  and the Hessian matrix  $\mathbf{H} \tilde{f}(\tilde{\mu})$  is non-negative definite and symmetric. Moreover, since  $\tilde{f}$  is a polynomial splines of order  $q_k \geq \alpha_k > 2$ ,  $k = 1, \dots, d$ , by Assumption 2,  $\mathbf{H} \tilde{f}(\mathbf{x})$  is continuous. Now since  $\tilde{f} \rightarrow f_0$  uniformly in  $P_0$ -probability as a result of Theorem A.5 and  $f_0$  has a well-separated maximum

by Assumption 2, it follows that  $\tilde{\mu}$  is consistent in estimating  $\mu_0$  by Theorem 5.7 of van der Vaart [30]. Hence for  $n$  large enough,  $\mathcal{B}(\tilde{\mu}, \tau/2)$  is contained in  $\mathcal{B}(\mu_0, \tau)$  for the same  $\tau$  appearing in Assumption 3. Now since  $D^r \tilde{f}$  converges uniformly to  $D^r f_0$  (Theorem A.5), and using the fact that maximum eigenvalue is a continuous operation, we have by the continuous mapping theorem that  $\lambda_{\max}(\mathbf{H} \tilde{f}(\mathbf{x})) \rightarrow \lambda_{\max}(\mathbf{H} f_0(\mathbf{x}))$  uniformly in  $\mathbf{x}$ . By further adapting (4.2) to the present situation, we see that  $\sup_{\mathbf{x} \in \mathcal{B}(\tilde{\mu}, \tau/2)} \lambda_{\max}(\mathbf{H} \tilde{f}(\mathbf{x})) < -\lambda_0$  for sufficiently large  $n$ . Consequently, the proof of the inequality (4.1) given in Theorem 4.1 will go through even if we replace  $f_0$  with  $\tilde{f}$  and  $\mu_0$  with  $\tilde{\mu}$ .

Let  $\mathbf{1}_d$  be a  $d$ -dimensional vector of ones, and  $\xi$  some point between  $\mathbf{x}$  and  $\mathbf{x} - (Rd)^{-1} \times \mathbf{1}_d \rho R_{n,k,\gamma}$  for any given  $\mathbf{x} \in [0, 1]^d$ . In what follows, we take  $n$  large enough so that  $R_{n,k,\gamma}$  is small enough, and adding or subtracting some constant multiple of  $R_{n,k,\gamma}$  still allows us to stay within  $[0, 1]^d$ . For (iii), we have by the multivariate mean value theorem and the Cauchy–Schwarz inequality that

$$|D_k \tilde{f}(\mathbf{x} - (Rd)^{-1} \mathbf{1}_d \rho R_{n,k,\gamma}) - D_k \tilde{f}(\mathbf{x})| \leq \|\nabla D_k \tilde{f}(\xi)\| R^{-1} d^{-1/2} \rho R_{n,k,\gamma}$$

for  $k = 1, \dots, d$ . Since  $D^r \tilde{f} \rightarrow D^r f_0$  uniformly (cf. Theorem A.5) and the norm is a continuous map,  $\|\nabla D_k \tilde{f}(\xi)\| \rightarrow \|\nabla D_k f_0(\xi)\| \leq \sqrt{d} \max_{1 \leq j \leq d} |D_j D_k f_0(\xi)| \leq \sqrt{d} \|f_0\|_{\alpha, \infty}$  in view of the definition given in (2.1). Therefore uniformly over  $\|f_0\|_{\alpha, \infty} \leq R$ , the right hand side above is less than  $\rho R_{n,k,\gamma}$  when  $n$  is large enough. Now since the mode of  $\tilde{f}(\cdot - (Rd)^{-1} \mathbf{1}_d \rho R_{n,k,\gamma})$  is  $\tilde{\mu} + (Rd)^{-1} \mathbf{1}_d \rho R_{n,k,\gamma}$ , it follows immediately that  $\tilde{\mu} + (Rd)^{-1} \mathbf{1}_d \rho \max_{1 \leq k \leq d} R_{n,k,\gamma} \in \mathcal{C}_\mu$ . Notice that this argument still holds even when we replace  $\mathbf{1}_d$  with any point in the boundary of a unit cube, and this collectively shows that with probability tending to one,  $\mathcal{C}_\mu$  contains a hyper-cube centered at  $\tilde{\mu}$  of size  $(Rd)^{-1} \rho \max_{1 \leq k \leq d} R_{n,k,\gamma}$ .

The proof of (iv) is similar to that of (i) with  $\max_{1 \leq k \leq d} R_{n,k,\gamma}$  replaced by  $R_{n,0,\gamma}$ . The proof of assertion (v) can be completed by following the arguments used in the proof of assertion (ii) with (4.2) applied to the pair  $f$  and  $\tilde{f}$ .  $\square$

To prove the results in Section 5, we need to first lay out some preliminary details. Define  $f_{0,z}(\mathbf{x} - \tilde{\mu}) = f_0(\mathbf{x})$  to be the shifted true function. Let  $\theta_0 = (\theta_{0,i} : i \leq m_\alpha)^T$  be a random vector such that  $f_{\theta_0}(\mathbf{x} - \tilde{\mu}) = T_{\mu_0} f_0(\mathbf{x})$ , where  $f_\theta$  is from (5.2) and  $T_{\mu_0} f_0(\mathbf{x})$  is the Taylor polynomial of order  $m_\alpha$  by expanding  $f_0$  around  $\mu_0$ , that is,

$$\sum_{i \leq m_\alpha} \theta_{0,i} (\mathbf{x} - \tilde{\mu})^i = f_0(\mu_0) + \sum_{i \leq m_\alpha, |i| \geq 2} \frac{D^i f_0(\mu_0)}{i!} (\mathbf{x} - \mu_0)^i, \quad (8.2)$$

where  $\nabla f_0(\mu_0) = \mathbf{0}$  by Assumption 2. Hence,  $\theta_0$  can be thought of as the true  $\theta$  by projecting  $f_0$  onto the space of polynomials of order  $m_\alpha$ . Note that  $\theta_0$  is random and depends on  $\tilde{\mu}$ ,  $\mu_0$  and  $f_0$ . By applying  $D^i$  on both sides of (8.2) and evaluating at  $\mathbf{x} = \tilde{\mu}$ , we have  $i! \theta_{0,i} = D^i T_{\mu_0} f_0(\tilde{\mu})$ . Note that since  $D^i f_0(\mathbf{x})$ ,  $i \leq m_\alpha$ , are continuous by Assumption 2, they are bounded over  $\{\mu : |\mu_k - \tilde{\mu}_k| \leq \delta_{n,k}, k = 1, \dots, d\}$ , and this implies that for any  $i \leq m_\alpha$ ,  $|\theta_{0,i}| = O_{P_0}(1)$  uniformly over  $\|f_0\|_{\alpha, \infty} \leq R$ . The design matrix  $\mathbf{Z}$  is generated using i.i.d. uniform samples and by Lemma 8.1, we know that  $\mathbf{Z}^T \mathbf{Z}$  is invertible with probability going to 1 as  $n \rightarrow \infty$ . In the actual computation, the invertibility of  $\mathbf{Z}^T \mathbf{Z}$  is not an important issue as there is the presence of the

prior covariance matrix  $\mathbf{V}$  to serve as a regularization factor, and  $\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1}$  in (5.3) is always invertible by our choice of  $\mathbf{V}$ .

As a consequence of Theorem A.5,  $D_k \tilde{f}$  converges uniformly to  $D_k f_0$  at the rate  $\epsilon_{n,k}$ , and by (8.1), this translates to  $\|\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}_0\| = O_{P_0}(\epsilon_n)$ . Let  $\mathbf{F}_0 = (f_0(\mathbf{x}_1), \dots, f_0(\mathbf{x}_{n_2}))^T$  where  $\{\mathbf{x}_1, \dots, \mathbf{x}_{n_2}\}$  is the original (unshifted) second stage samples. Note that  $(\mathbf{Z}\boldsymbol{\theta}_0)_i = f_{\theta_0}(\mathbf{x}_i - \tilde{\boldsymbol{\mu}}) = T_{\boldsymbol{\mu}_0} f_0(\mathbf{x}_i)$  and under the assumption of (2.2), we have

$$\begin{aligned} \|\mathbf{F}_0 - \mathbf{Z}\boldsymbol{\theta}_0\|_\infty &= \max_{1 \leq i \leq n_2} |f_0(\mathbf{x}_i) - T_{\boldsymbol{\mu}_0} f_0(\mathbf{x}_i)| \lesssim \max_{1 \leq i \leq n_2} \sum_{k=1}^d |x_{ik} - \mu_{0k}|^{\alpha_k} \\ &\lesssim \max_{1 \leq i \leq n_2} \sum_{k=1}^d |x_{ik} - \tilde{\mu}_k|^{\alpha_k} + \sum_{k=1}^d |\tilde{\mu}_k - \mu_{0k}|^{\alpha_k} \lesssim \sum_{k=1}^d \delta_{n,k}^{\alpha_k} \end{aligned} \quad (8.3)$$

uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$ . The second line follows from the inequality  $|x + y|^r \leq \max\{1, 2^{r-1}\}(|x|^r + |y|^r)$ , while the last inequality is due to  $|x_{ik} - \tilde{\mu}_k| \leq \delta_{n,k}$  almost surely since  $x_{ik} - \tilde{\mu}_k \sim \text{Uniform}(-\delta_{n,k}, \delta_{n,k})$ ; and  $|\tilde{\mu}_k - \mu_{0k}| \leq \|\tilde{\boldsymbol{\mu}} - \boldsymbol{\mu}_0\| = O_{P_0}(\epsilon_n) = o_{P_0}(\delta_{n,k})$  as argued above, where  $\epsilon_n = o(\delta_{n,k})$ ,  $k = 1, \dots, d$  was by our choice in (5.1).

We break the proof of Theorem 5.2 into a series of steps. First, let us enumerate the elements of  $\{i : i \leq m_\alpha\}$  as  $\{i_0, \dots, i_W\}$  with  $W + 1 = \prod_{k=1}^d \alpha_k$ . For the rest of this section, we follow this indexing convention and index the entries of vectors  $\boldsymbol{\xi}$ ,  $\boldsymbol{\theta}$  and  $\boldsymbol{\theta}_0$  by elements of  $\{i : i \leq m_\alpha\}$ . For matrices  $\mathbf{Z}^T \mathbf{Z}$  and  $\mathbf{V}$ , we enumerate their rows and columns starting from 0 and ending at  $W$ . We note that in our first and second stage sampling plans,  $n_1 \asymp n \asymp n_2$ .

The first key step is to derive sharp upper bounds for the posterior mean and variance, which will involve upper bounding the entries of  $(\mathbf{Z}^T \mathbf{Z})^{-1}$ . These calculations are made simpler by centering the design points so that they are uniformly distributed around zero in each co-ordinate, and  $(\mathbf{Z}^T \mathbf{Z})^{-1}$  will not depend on  $\tilde{\boldsymbol{\mu}}$ . The following lemma describes the asymptotic behavior of the entries of  $(\mathbf{Z}^T \mathbf{Z})^{-1}$  when the second stage samples are collected under uniform random sampling. In the following, we define  $\boldsymbol{\delta}_n = (\delta_{n,1}, \dots, \delta_{n,d})^T$  and write  $\boldsymbol{\delta}_n^i$  to mean  $\prod_{k=1}^d \delta_{n,k}^{i_k}$ .

**Lemma 8.1.** *As  $n \rightarrow \infty$ ,  $\mathbf{Z}^T \mathbf{Z}$  is invertible with probability tending to 1. Moreover for  $a, b = 0, \dots, W$ , we have  $[(\mathbf{Z}^T \mathbf{Z})^{-1}]_{ab} = O_P(n^{-1} \boldsymbol{\delta}_n^{-(i_a + i_b)})$ .*

**Proof.** After centering, second stage samples are distributed as  $\mathbf{z}_i = (z_{i1}, \dots, z_{id})^T \stackrel{\text{i.i.d.}}{\sim} \text{Uniform}(\prod_{k=1}^d [-\delta_{n,k}, \delta_{n,k}])$ ,  $i = 1, \dots, n_2$ , and  $z_{ik} \sim \text{Uniform}[-\delta_{n,k}, \delta_{n,k}]$ . Thus,

$$\mathbf{Z}^T \mathbf{Z} = n_2 \begin{pmatrix} a_{00} & a_{01} \boldsymbol{\delta}_n^{i_1} & \cdots & a_{0W} \boldsymbol{\delta}_n^{i_W} \\ a_{10} \boldsymbol{\delta}_n^{i_1} & a_{11} \boldsymbol{\delta}_n^{2i_1} & \cdots & a_{1W} \boldsymbol{\delta}_n^{i_1 + i_W} \\ \vdots & \vdots & \ddots & \vdots \\ a_{W0} \boldsymbol{\delta}_n^{i_W} & a_{W1} \boldsymbol{\delta}_n^{i_W + i_1} & \cdots & a_{WW} \boldsymbol{\delta}_n^{2i_W} \end{pmatrix} := n_2 \boldsymbol{\Delta} \mathbf{A} \boldsymbol{\Delta} \quad (8.4)$$

with the  $(i, j)$ -entry of  $\mathbf{A}$  being  $a_{ij} = n_2^{-1} \sum_{k=1}^{n_2} \mathbf{U}_k^{i_i} \mathbf{U}_k^{i_j}$  where  $\mathbf{U}_k = (U_{k1}, \dots, U_{kd})^T \stackrel{\text{i.i.d.}}{\sim} \text{Uniform}[-1, 1]^d$ , and  $\boldsymbol{\Delta} = \text{diag}\{\boldsymbol{\delta}_n^{i_j} : j = 0, \dots, W\}$ . Define  $\mathbb{U} = (\mathbf{U}^{i_0}, \dots, \mathbf{U}^{i_W})^T$  for  $\mathbf{U} =$



$(U_1, \dots, U_d)^T \sim \text{Uniform}[-1, 1]^d$ . By the law of large numbers, we have that  $\mathbf{A}$  converges in probability to  $\mathbf{E}\mathbf{U}\mathbf{U}^T$  entry-wise, and hence  $\mathbf{E}\mathbf{U}\mathbf{U}^T - \epsilon \mathbf{I} \leq \mathbf{A} \leq \mathbf{E}\mathbf{U}\mathbf{U}^T + \epsilon \mathbf{I}$  for a sufficiently small  $\epsilon > 0$ . Observe that entries of  $\mathbf{E}\mathbf{U}\mathbf{U}^T$  are mixed moments of  $U \sim \text{Uniform}[-1, 1]$  and hence is positive definite. Then  $\mathbf{E}\mathbf{U}\mathbf{U}^T - \epsilon \mathbf{I}$  is invertible when  $\epsilon$  is smaller than the minimum eigenvalue of  $\mathbf{E}\mathbf{U}\mathbf{U}^T$ . Thus for sufficiently small  $\epsilon > 0$ ,

$$n_2^{-1} \Delta^{-1} (\mathbf{E}\mathbf{U}\mathbf{U}^T + \epsilon \mathbf{I})^{-1} \Delta^{-1} \leq (\mathbf{Z}^T \mathbf{Z})^{-1} \leq n_2^{-1} \Delta^{-1} (\mathbf{E}\mathbf{U}\mathbf{U}^T - \epsilon \mathbf{I})^{-1} \Delta^{-1}.$$

Let  $u^{ij}$  be the  $(i, j)$ th entry of  $(\mathbf{E}\mathbf{U}\mathbf{U}^T)^{-1}$  and recall that  $n_2 \geq cn$  for some constant  $c > 0$ . It then follows that for  $a = 0, \dots, W$ ,  $[(\mathbf{Z}^T \mathbf{Z})^{-1}]_{aa}$  is  $O_P(n_2^{-1} u^{aa} \delta_n^{-(i_a+i_a)}) = O_P(n^{-1} \delta_n^{-(i_a+i_a)})$ . Using the fact that for positive definite  $\mathbf{G}$ ,  $g_{ij} \leq \sqrt{g_{ii} g_{jj}}$  by the Cauchy-Schwarz inequality,  $[(\mathbf{Z}^T \mathbf{Z})^{-1}]_{ab}$  with  $a, b = 0, \dots, W$  is bounded above by

$$\sqrt{[(\mathbf{Z}^T \mathbf{Z})^{-1}]_{aa} [(\mathbf{Z}^T \mathbf{Z})^{-1}]_{bb}} = O_P(n^{-1} \delta_n^{-(i_a+i_b)}). \quad \square$$

**Proof of Proposition 5.1.** See supplementary article Yoo and Ghosal [34].  $\square$

Let  $\mathcal{K}_n := [\sigma_0^2 - m_n \xi_n, \sigma_0^2 + m_n \xi_n]$  where  $m_n$  is any sequence going to infinity (e.g., slowly varying such as  $\log n$ ) and  $\xi_n = \max\{n^{-1/2}, n^{-2\alpha^*/(2\alpha^*+d)}, \sum_{k=1}^d \delta_{n,k}^{2\alpha_k}\}$ . Recall that  $\mathcal{Q} = \{\mathbf{x} : |x_k| \leq \delta_{n,k}, k = 1, \dots, d\}$  is the centered second stage credible set/sampling region.

**Lemma 8.2.** For any  $\mathbf{i} \leq \mathbf{m}_\alpha$  and  $\sigma^2 \in \mathcal{K}_n$ ,

$$\mathbb{E}[(\theta_i - \theta_{0,i})^2 | \mathbf{Y}, \sigma^2] = O_{P_0} \left[ \prod_{k=1}^d \delta_{n,k}^{-2i_k} \left( \frac{1}{n} + \sum_{k=1}^d \delta_{n,k}^{2\alpha_k} \right) \right].$$

**Proof.** Let  $0 \leq h \leq W$ . Since  $\mathbf{V} > \mathbf{0}$  by assumption, we have by Lemma 8.1 that  $\sup_{\sigma^2 \in \mathcal{K}_n} \text{Var}(\theta_{ih} | \mathbf{Y}, \sigma^2)$  is

$$[\sigma_0^2 + o(1)] [(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}]_{hh} \lesssim [(\mathbf{Z}^T \mathbf{Z})^{-1}]_{hh} \lesssim n^{-1} \delta_n^{-2i_h}. \quad (8.5)$$

Now the bias for the conditional posterior mean in vector form is

$$\begin{aligned} \mathbb{E}(\boldsymbol{\theta} | \mathbf{Y}, \sigma^2) - \boldsymbol{\theta}_0 &= (\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1} (\mathbf{Z}^T \mathbf{Y} + \mathbf{V}^{-1} \boldsymbol{\xi}) - \boldsymbol{\theta}_0 \\ &= (\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1} [\mathbf{Z}^T \boldsymbol{\varepsilon} + \mathbf{Z}^T (\mathbf{F}_0 - \mathbf{Z} \boldsymbol{\theta}_0) + \mathbf{V}^{-1} (\boldsymbol{\xi} - \boldsymbol{\theta}_0)]. \end{aligned} \quad (8.6)$$

Following the same reasoning as in (8.5), the  $h$ th diagonal entry of the covariance matrix  $\sigma_0^2 (\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1} \mathbf{Z}^T \mathbf{Z} (\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}$  of  $(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1} \boldsymbol{\varepsilon}$  is

$$\sigma_0^2 [(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1} \mathbf{Z}^T \mathbf{Z} (\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}]_{hh} \lesssim n^{-1} \delta_n^{-2i_h}.$$

Since  $\mathbb{E}_0(\boldsymbol{\varepsilon}) = \mathbf{0}$ , it follows from Markov's inequality that the  $h$ th entry of  $(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1} \mathbf{Z}^T \boldsymbol{\varepsilon}$  is  $O_{P_0}(n^{-1/2} \delta_n^{-i_h})$  for  $0 \leq h \leq W$ .

Let  $\beta_{ij}$  be the  $(i, j)$ th element of  $(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}$ ,  $\kappa_{ij}$  be the  $(i, j)$ th element of  $\mathbf{V}^{-1}$  and  $\gamma_i$  be the  $i$ th entry of  $\mathbf{F}_0 - \mathbf{Z}\boldsymbol{\theta}_0$ . By (8.3), we have uniformly over  $1 \leq i \leq n$  that  $|\gamma_i| \lesssim \sum_{k=1}^d \delta_{n,k}^{\alpha_k}$ . Now using the fact that for positive definite  $\mathbf{G}$ ,  $g_{hj} \leq \sqrt{g_{hh}g_{jj}}$  by the Cauchy–Schwarz inequality, we have for  $0 \leq h, j \leq W$ ,

$$\begin{aligned} [(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}]_{hj} &\leq \sqrt{[(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}]_{hh}[(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1}]_{jj}} \\ &\leq \sqrt{[(\mathbf{Z}^T \mathbf{Z})^{-1}]_{hh}[(\mathbf{Z}^T \mathbf{Z})^{-1}]_{jj}} \lesssim n^{-1} \delta_n^{-(i_h+i_j)}. \end{aligned}$$

Since  $\mathbf{z}_j \in \mathcal{Q}$ ,  $j = 1, \dots, n_2$ , we have  $|\mathbf{z}_j^i| \leq \delta_n^i$  for  $\mathbf{i} \leq \mathbf{m}_\alpha$ . Therefore, since  $n_2 \leq n$ ,  $[(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1} \mathbf{Z}^T (\mathbf{F}_0 - \mathbf{Z}\boldsymbol{\theta}_0)]_h$  is

$$\beta_{h0} \sum_{j=1}^{n_2} \mathbf{z}_j^{i_0} \gamma_j + \dots + \beta_{hW} \sum_{j=1}^{n_2} \mathbf{z}_j^{i_W} \gamma_j \lesssim \delta_n^{-i_h} \sum_{k=1}^d \delta_{n,k}^{\alpha_k}.$$

It remains to bound each entry of the last term in (8.6). Since  $|\theta_{0,i_j} - \xi_{0,i_j}| \leq |\theta_{0,i_j}| + |\xi_{0,i_j}| = O_{P_0}(1)$  for  $j = 0, \dots, W$ , then Lemma 8.1 and the choice of  $\mathbf{V}$  imply that  $[(\mathbf{Z}^T \mathbf{Z} + \mathbf{V}^{-1})^{-1} \times \mathbf{V}^{-1}(\boldsymbol{\xi} - \boldsymbol{\theta}_0)]_h$  is

$$\beta_{h0} \sum_{j=0}^W \kappa_{0j}(\xi_{i_j} - \theta_{0,i_j}) + \dots + \beta_{hW} \sum_{j=0}^W \kappa_{Wj}(\xi_{i_j} - \theta_{0,i_j}) \lesssim n^{-1} \delta_n^{-i_h}.$$

Combining the bounds derived back into (8.6), the squared bias  $[E(\theta_{i_h}|\mathbf{Y}, \sigma^2) - \theta_{0,i_h}]^2$  is  $O_{P_0}[n^{-2} \delta_n^{-2i_h} + \delta_n^{-2i_h} \sum_{k=1}^d \delta_{n,k}^{2\alpha_k}]$ . The result follows in view of the bounds established and (8.5).  $\square$

**Lemma 8.3.** *Uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$ , for any  $\mathbf{r} \leq \mathbf{m}_\alpha$ ,  $\mathbf{x} \in \mathcal{Q}$  and  $m_n \rightarrow \infty$ ,  $E_0 \Pi(|D^{\mathbf{r}} f_{\boldsymbol{\theta}}(\mathbf{x}) - D^{\mathbf{r}} f_{\boldsymbol{\theta}_0}(\mathbf{x})| > m_n \epsilon_{n,\mathbf{r}} | \mathbf{Y}) \rightarrow 0$ , where  $\epsilon_{n,\mathbf{r}} := \delta_n^{-\mathbf{r}} (n^{-1/2} + \sum_{k=1}^d \delta_{n,k}^{\alpha_k})$ .*

**Proof.** In view of (5.4),

$$D^{\mathbf{r}} f_{\boldsymbol{\theta}}(\mathbf{x}) - D^{\mathbf{r}} f_{\boldsymbol{\theta}_0}(\mathbf{x}) = \mathbf{r}!(\boldsymbol{\theta}_{\mathbf{r}} - \boldsymbol{\theta}_{0,\mathbf{r}}) + \sum_{\mathbf{r} \leq \mathbf{i} \leq \mathbf{m}_\alpha, \mathbf{i} \neq \mathbf{r}} \frac{\mathbf{i}!}{(\mathbf{i} - \mathbf{r})!} (\boldsymbol{\theta}_{\mathbf{i}} - \boldsymbol{\theta}_{0,\mathbf{i}}) \mathbf{x}^{\mathbf{i}-\mathbf{r}}.$$

Observe that for any  $\mathbf{x} \in \mathcal{Q}$ ,  $|\mathbf{x}^{\mathbf{i}-\mathbf{r}}| \leq \delta_n^{i-\mathbf{r}}$ . Also, by noting that  $r_k \leq i_k \leq \alpha_k - 1$  for  $k = 1, \dots, d$ , we have both  $\mathbf{r}!, \mathbf{i}! \leq \prod_{k=1}^d (\alpha_k - 1)$ . Using the fact  $(\sum_{i=1}^n |b_i|)^p \leq n^{p-1} \sum_{i=1}^n |b_i|^p$  for  $p \geq 1$ ,  $|D^{\mathbf{r}} f_{\boldsymbol{\theta}}(\mathbf{x}) - D^{\mathbf{r}} f_{\boldsymbol{\theta}_0}(\mathbf{x})|^2$  is bounded above up to a constant multiple by

$$|\boldsymbol{\theta}_{\mathbf{r}} - \boldsymbol{\theta}_{0,\mathbf{r}}|^2 + \sum_{\mathbf{r} \leq \mathbf{i} \leq \mathbf{m}_\alpha, \mathbf{i} \neq \mathbf{r}} |\boldsymbol{\theta}_{\mathbf{i}} - \boldsymbol{\theta}_{0,\mathbf{i}}|^2 \delta_n^{2i-2\mathbf{r}}. \quad (8.7)$$

Therefore, for any  $\mathbf{r} \leq \mathbf{m}_\alpha$  and any  $\mathbf{x} \in \mathcal{Q}$ , we have uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$  that  $E_0 \sup_{\sigma^2 \in \mathcal{K}_n} E(|D^{\mathbf{r}} f_\theta(\mathbf{x}) - D^{\mathbf{r}} f_{\theta_0}(\mathbf{x})|^2 | \mathbf{Y}, \sigma^2)$  is bounded up to a constant multiple by

$$E_0 \sup_{\sigma^2 \in \mathcal{K}_n} E[(\theta_{\mathbf{r}} - \theta_{0,\mathbf{r}})^2 | \mathbf{Y}, \sigma^2] + \sum_{\mathbf{r} \leq \mathbf{i} \leq \mathbf{m}_\alpha, \mathbf{i} \neq \mathbf{r}} \delta_n^{2\mathbf{i}-2\mathbf{r}} E_0 \sup_{\sigma^2 \in \mathcal{K}_n} E[(\theta_{\mathbf{i}} - \theta_{0,\mathbf{i}})^2 | \mathbf{Y}, \sigma^2]. \quad (8.8)$$

In view of Lemma 8.2, the first term is bounded by  $\delta_n^{-2\mathbf{r}}(n^{-1} + \sum_{k=1}^d \delta_{n,k}^{2\alpha_k})$ . By noting that the sum over  $\{\mathbf{i} : \mathbf{r} \leq \mathbf{i} \leq \mathbf{m}_\alpha, \mathbf{i} \neq \mathbf{r}\}$  has at most  $\prod_{k=1}^d \alpha_k$  terms, another application of Lemma 8.2 implies that the second term in (8.8) is bounded above by

$$\sum_{\mathbf{r} \leq \mathbf{i} \leq \mathbf{m}_\alpha, \mathbf{i} \neq \mathbf{r}} \delta_n^{2\mathbf{i}-2\mathbf{r}} \left[ \delta_n^{-2\mathbf{i}} \left( \frac{1}{n} + \sum_{k=1}^d \delta_{n,k}^{2\alpha_k} \right) \right] \lesssim \delta_n^{-2\mathbf{r}} \left( \frac{1}{n} + \sum_{k=1}^d \delta_{n,k}^{2\alpha_k} \right).$$

Using the inequality  $|a+b|^r \leq \max(1, 2^{r-1})(|a|^r + |b|^r)$ ,  $r > 0$  and (2.2), we have  $|D^{\mathbf{r}} f_{\theta_0}(\mathbf{x}) - D^{\mathbf{r}} f_{0,z}(\mathbf{x})|$  is

$$\begin{aligned} |D^{\mathbf{r}} T_{\mu_0}(\mathbf{x} + \tilde{\boldsymbol{\mu}}) - D^{\mathbf{r}} f_0(\mathbf{x} + \tilde{\boldsymbol{\mu}})| &\lesssim \sum_{k=1}^d |x_k + \tilde{\mu}_k - \mu_{0,k}|^{\alpha_k - r_k} \\ &\lesssim \sum_{k=1}^d \delta_{n,k}^{\alpha_k - r_k} + \sum_{k=1}^d |\tilde{\mu}_k - \mu_{0,k}|^{\alpha_k - r_k}, \end{aligned} \quad (8.9)$$

and  $E_0 |D^{\mathbf{r}} f_{\theta_0}(\mathbf{x}) - D^{\mathbf{r}} f_{0,z}(\mathbf{x})|^2 \lesssim \sum_{k=1}^d \delta_{n,k}^{2\alpha_k - 2r_k}$  uniformly in  $\|f_0\|_{\alpha,\infty} \leq R$  by (8.3).

Define  $P_{n,r}(\mathbf{x}) := E_0 \sup_{\sigma^2 \in \mathcal{K}_n} E[(D^{\mathbf{r}} f_\theta(\mathbf{x}) - D^{\mathbf{r}} f_{0,z}(\mathbf{x}))^2 | \mathbf{Y}, \sigma^2]$ . Combining all the bounds established and (8.8), we have uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$ ,

$$\begin{aligned} P_{n,r}(\mathbf{x}) &\lesssim E_0 \sup_{\sigma^2 \in \mathcal{K}_n} E(|D^{\mathbf{r}} f_\theta(\mathbf{x}) - D^{\mathbf{r}} f_{\theta_0}(\mathbf{x})|^2 | \mathbf{Y}, \sigma^2) + E_0 |D^{\mathbf{r}} f_{\theta_0}(\mathbf{x}) - D^{\mathbf{r}} f_{0,z}(\mathbf{x})|^2 \\ &\lesssim \delta_n^{-2\mathbf{r}} \left( \frac{1}{n} + \sum_{k=1}^d \delta_{n,k}^{2\alpha_k} \right) + \sum_{k=1}^d \delta_{n,k}^{2\alpha_k - 2r_k} \lesssim \epsilon_{n,r}^2. \end{aligned}$$

By Proposition 5.1,  $P_0(\tilde{\sigma}_*^2 \in \mathcal{K}_n) \rightarrow 1$  as  $n \rightarrow \infty$  uniformly in  $\|f_0\|_{\alpha,\infty} \leq R$ . For the empirical Bayes posterior  $\Pi(\cdot | \mathbf{Y}) \equiv \Pi_{\tilde{\sigma}_*}(\cdot | \mathbf{Y})$  and by Markov's inequality, we have for any  $m_n \rightarrow \infty$ ,

$$E_0 \Pi_{\tilde{\sigma}_*}(|D^{\mathbf{r}} f_\theta(\mathbf{x}) - D^{\mathbf{r}} f_{0,z}(\mathbf{x})| > m_n \epsilon_{n,r} | \mathbf{Y}) \leq \frac{P_{n,r}(\mathbf{x})}{m_n^2 \epsilon_{n,r}^2} + o(1) \rightarrow 0, \quad (8.10)$$

uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$ . For the hierarchical Bayes procedure, we have for any  $m_n \rightarrow \infty$  that  $E_0 \Pi(|D^{\mathbf{r}} f_\theta(\mathbf{x}) - D^{\mathbf{r}} f_{0,z}(\mathbf{x})| > m_n \epsilon_{n,r} | \mathbf{Y})$  is uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$  bounded above by

$$P_{n,r}(\mathbf{x}) / (m_n \epsilon_{n,r})^2 + E_0 \Pi(\sigma^2 \notin \mathcal{K}_n | \mathbf{Y}). \quad (8.11)$$

The first term is  $o(1)$  since  $P_{n,r}(\mathbf{x}) \lesssim \epsilon_{n,r}^2$ , while the second term goes to zero by Proposition 5.1.  $\square$

An immediate consequence of the previous lemma is the following, whose proof can be found in the supplementary article Yoo and Ghosal [34].

**Corollary 8.4.** *Uniformly in  $\|f_0\|_{\alpha,\infty} \leq R$ , we have for any  $\mathbf{r} \leq \mathbf{m}_\alpha$  and  $m_n \rightarrow \infty$  that  $E_0 \Pi(\|D^{\mathbf{r}} f_\theta - D^{\mathbf{r}} f_{0,z}\|_\infty > m_n \epsilon_{n,r} | \mathbf{Y}) \rightarrow 0$ .*

We are now ready to prove Theorem 5.2.

**Proof of Theorem 5.2.** We shall prove only the empirical Bayes case as the hierarchical Bayes case follows the same steps. Recall that  $\boldsymbol{\mu} = \tilde{\boldsymbol{\mu}} + \boldsymbol{\mu}_z$ . As a consequence of Theorem 4.1 and our choice of  $\delta_{n,k}$ ,  $k = 1, \dots, d$ , we have  $P_0(\boldsymbol{\mu}_0 - \tilde{\boldsymbol{\mu}} \in \mathcal{Q}) \rightarrow 1$ . Therefore by (4.1),

$$\|\boldsymbol{\mu} - \boldsymbol{\mu}_0\| = \|\boldsymbol{\mu}_z - (\boldsymbol{\mu}_0 - \tilde{\boldsymbol{\mu}})\| \leq \frac{\sqrt{d}}{\lambda_0} \max_{1 \leq k \leq d} \sup_{\mathbf{x} \in \mathcal{Q}} |D_k f_\theta(\mathbf{x}) - D_k f_{0,z}(\mathbf{x})|.$$

Let  $\tau_{n,k} := \delta_{n,k}^{-1}(n^{-1/2} + \sum_{k=1}^d \delta_{n,k}^{\alpha_k})$ . Using this bound and Corollary 8.4 with  $\mathbf{r} = \mathbf{e}_k$ , we have for any  $m_n \rightarrow \infty$  that  $E_0 \Pi(\|\boldsymbol{\mu} - \boldsymbol{\mu}_0\| > m_n \max_{1 \leq k \leq d} \tau_{n,k} | \mathbf{Y})$  is bounded above by

$$\sum_{k=1}^d E_0 \Pi\left(\|D_k f_\theta - D_k f_{0,z}\|_\infty > \frac{\lambda_0}{\sqrt{d}} m_n \max_{1 \leq k \leq d} \tau_{n,k} \mid \mathbf{Y}\right) \rightarrow 0.$$

By definition,  $M = f_\theta(\boldsymbol{\mu}_z)$  and  $M_0 = f_{0,z}(\boldsymbol{\mu}_0 - \tilde{\boldsymbol{\mu}})$ . Then by (4.2),  $|M - M_0| \leq \sup_{\mathbf{x} \in \mathcal{Q}} |f_\theta(\mathbf{x}) - f_{0,z}(\mathbf{x})|$  since  $P_0(\boldsymbol{\mu}_0 - \tilde{\boldsymbol{\mu}} \in \mathcal{Q}) \rightarrow 1$  as before. Therefore by Corollary 8.4 with  $\mathbf{r} = \mathbf{0}$ , we have for  $m_n \rightarrow \infty$ ,  $E_0 \Pi[|M - M_0| > m_n(n^{-1/2} + \sum_{k=1}^d \delta_{n,k}^{\alpha_k}) | \mathbf{Y}] \leq E_0 \Pi[\|f_\theta - f_{0,z}\|_\infty > m_n(n^{-1/2} + \sum_{k=1}^d \delta_{n,k}^{\alpha_k}) | \mathbf{Y}] \rightarrow 0$ , uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$ .

To prove the last part, note that  $\delta_{n,k} = n^{-1/(2\alpha_k)}$ ,  $k = 1, \dots, d$ , comes from equating the two terms in the second stage rates for  $\boldsymbol{\mu}$  and  $M$ , that is,  $n^{-1/2} = \sum_{k=1}^d \delta_{n,k}^{\alpha_k}$ . To solve for  $\delta_{n,k}$ , take  $\delta_{n,k} = (d^{-1} \delta_n)^{1/\alpha_k}$ , where  $\delta_n$  is a positive sequence in  $n$  that does not depend on  $k$ . It follows that  $\delta_n = n^{-1/2}$  and hence  $\delta_{n,k} = n^{-1/(2\alpha_k)}$ . The condition  $\min_{1 \leq k \leq d} \delta_{n,k} = \rho_n \epsilon_n$  or equivalently  $\epsilon_n = o(\min_{1 \leq k \leq d} \delta_{n,k})$  is fulfilled when  $1/(2\underline{\alpha}) < \alpha^*(1 - \underline{\alpha}^{-1})/(2\alpha^* + d)$  for  $\underline{\alpha} = \min_{1 \leq k \leq d} \alpha_k$ . By rearranging, we need  $2\underline{\alpha}\alpha^* - 4\alpha^* - d > 0$ . Since  $\underline{\alpha} > 2$  by Assumption 2, we have  $2\underline{\alpha}\alpha^* - 4\alpha^* - d = 2(\underline{\alpha} - 2)\alpha^* - d \geq 2\underline{\alpha}^2 - 4\underline{\alpha} - d$ . Thus, it suffices to find  $\underline{\alpha}$  such that  $2\underline{\alpha}^2 - 4\underline{\alpha} - d > 0$  and this is satisfied if  $\underline{\alpha} > 1 + \sqrt{1 + d/2}$  under the constraint that  $\underline{\alpha} > 2$ .  $\square$

## Appendix

In this section, we collect some auxiliary results and technical lemmas that were used in several places to prove the main theorems in the previous section.

The next two results concern contraction rates and credible band coverage for the regression function  $f$  and its derivatives, and they correspond to Theorems 4.4 and 5.3 of Yoo and Ghosal [33], respectively.

**Theorem A.5.** *Let  $J_{n,k} \asymp (n/\log n)^{\alpha^*/\{\alpha_k(2\alpha^*+d)\}}$ ,  $k = 1, \dots, d$ . Then under Assumption 1, we have for any sequence  $m_n \rightarrow \infty$ ,*

$$\sup_{\|f_0\|_{\alpha,\infty} \leq R} E_0 \Pi(\|D^{\mathbf{r}} f - D^{\mathbf{r}} f_0\|_{\infty} > m_n (\log n/n)^{\alpha^*\{1-\sum_{k=1}^d (r_k/\alpha_k)\}/(2\alpha^*+d)} | \mathbf{Y}) \rightarrow 0.$$

*In particular, contraction rate for  $f$  can be recovered by setting  $\mathbf{r} = \mathbf{0}$ .*

Consider the simultaneous credible band  $\{f : \|D^{\mathbf{r}} f - D^{\mathbf{r}} \tilde{f}\|_{\infty} \leq \rho R_{n,\mathbf{r},\gamma}\}$ , where the quantile  $R_{n,\mathbf{r},\gamma}$  is chosen such that  $\Pi(\|D^{\mathbf{r}} f - D^{\mathbf{r}} \tilde{f}\|_{\infty} \leq R_{n,\mathbf{r},\gamma} | \mathbf{Y}) = 1 - \gamma$  and  $\rho > 0$  is some large enough constant.

**Theorem A.6.** *Let  $J_{n,k} \asymp (n/\log n)^{\alpha^*/\{\alpha_k(2\alpha^*+d)\}}$ ,  $k = 1, \dots, d$ . Then under Assumption 1,*

1.  $\inf_{\|f_0\|_{\alpha,\infty} \leq R} P_0(\|D^{\mathbf{r}} \tilde{f} - D^{\mathbf{r}} f_0\|_{\infty} \leq \rho R_{n,\mathbf{r},\gamma}) \rightarrow 1$ ,
2.  $R_{n,\mathbf{r},\gamma} \asymp (\log n/n)^{\alpha^*\{1-\sum_{k=1}^d (r_k/\alpha_k)\}/(2\alpha^*+d)}$  in  $P_0$ -probability.

Credible bands for  $f$  is recovered by  $\mathbf{r} = \mathbf{0}$  and for  $D^{e_k} f \equiv D_k f$  by  $\mathbf{r} = \mathbf{e}_k$ , in the latter case we also write the radius as  $R_{n,k,\gamma} = R_{n,\mathbf{e}_k,\gamma}$ .

The result below was taken from (3.10) and (3.11) of Yoo and Ghosal [33], and it shows that B-splines despite being a non-orthonormal basis, are approximately orthogonal under our assumption on the design points.

**Lemma A.7.** *Let  $J = \prod_{k=1}^d J_k$ . If the design points were chosen such that (2.3) holds, then for some constants  $C_1, C_2, c_1, c_2 > 0$ ,*

$$C_1(n/J) \leq \mathbf{B}^T \mathbf{B} \leq C_2(n/J), \quad (\text{A.12})$$

$$C_1(n/J) + c_2^{-1} \leq \lambda_{\min}(\mathbf{B}^T \mathbf{B} + \mathbf{\Omega}^{-1}) \leq \lambda_{\max}(\mathbf{B}^T \mathbf{B} + \mathbf{\Omega}^{-1}) \leq C_2(n/J) + c_1^{-1}. \quad (\text{A.13})$$

**Lemma A.8.** *Let  $\mathbf{W}_{\mathbf{r}}$  be the finite difference matrix as seen in (3.1). Under the quasi-uniformity of the knot distribution, we have  $\lambda_{\max}(\mathbf{W}_{\mathbf{r}} \mathbf{W}_{\mathbf{r}}^T) = O(\prod_{k=1}^d J_k^{2r_k})$ .*

**Proof.** Note that  $\lambda_{\max}(\mathbf{W}_{\mathbf{r}} \mathbf{W}_{\mathbf{r}}^T) = \lambda_{\max}(\mathbf{W}_{\mathbf{r}}^T \mathbf{W}_{\mathbf{r}}) \leq \|\mathbf{W}_{\mathbf{r}}^T \mathbf{W}_{\mathbf{r}}\|_{(2,2)} \lesssim \prod_{k=1}^d J_k^{2r_k}$ , where the last upper bound was computed in (7.15) of Yoo and Ghosal [33].  $\square$

The following proposition shows that the single stage or the first stage (in the setting of two-stage procedure) empirical or hierarchical Bayes estimator of  $\sigma^2$  is consistent. Note that this result corresponds to Proposition 4.1 of Yoo and Ghosal [33].

**Proposition A.9 (First stage error variance).** *Suppose  $J_{n,k} \asymp (n/\log n)^{\alpha^*/\{\alpha_k(2\alpha^*+d)\}}$ . Then uniformly over  $\|f_0\|_{\alpha,\infty} \leq R$ ,*

- (a) First stage empirical Bayes estimator  $\tilde{\sigma}_1^2$  converges to  $\sigma_0^2$  under  $P_0$ -probability at the rate  $\max\{n^{-1/2}, n^{-2\alpha^*/(2\alpha^*+d)}\}$ .
- (b) If inverse gamma prior is used, first stage posterior for  $\sigma^2$  contracts to  $\sigma_0^2$  at the same rate.
- (c) If the prior used has continuous and positive density on  $(0, \infty)$ , then the first stage posterior for  $\sigma$  is consistent.

## Acknowledgement

We would like to thank the associate editor and the referees for their thought-provoking comments, as this led us to re-examine our results and improve our presentation.

## Supplementary Material

**Supplement to “Bayesian mode and maximum estimation and accelerated rates of contraction”** (DOI: [10.3150/18-BEJ1056SUPP](https://doi.org/10.3150/18-BEJ1056SUPP); .pdf). The supplementary file contains detailed proofs of Corollary 4.2, Proposition 5.1 and Corollary 8.4.

## References

- [1] Belitser, E., Ghosal, S. and van Zanten, H. (2012). Optimal two-stage procedures for estimating location and size of the maximum of a multivariate regression function. *Ann. Statist.* **40** 2850–2876. [MR3097962](#)
- [2] Blum, J.R. (1954). Multidimensional stochastic approximation methods. *Ann. Math. Statist.* **25** 737–744. [MR0065092](#)
- [3] Castillo, I. and Nickl, R. (2013). Nonparametric Bernstein–von Mises theorems in Gaussian white noise. *Ann. Statist.* **41** 1999–2028. [MR3127856](#)
- [4] Castillo, I. and Nickl, R. (2014). On the Bernstein–von Mises phenomenon for nonparametric Bayes procedures. *Ann. Statist.* **42** 1941–1969. [MR3262473](#)
- [5] Chen, H. (1988). Lower rate of convergence for locating a maximum of a function. *Ann. Statist.* **16** 1330–1334. [MR0959206](#)
- [6] Cox, D.D. (1993). An analysis of Bayesian inference for nonparametric regression. *Ann. Statist.* **21** 903–923. [MR1232525](#)
- [7] de Boor, C. (2001). *A Practical Guide to Splines*, Revised ed. New York: Springer. [MR1900298](#)
- [8] Dippon, J. (2003). Accelerated randomized stochastic optimization. *Ann. Statist.* **31** 1260–1281. [MR2001650](#)
- [9] Fabian, V. (1967). Stochastic approximation of minima with improved asymptotic speed. *Ann. Math. Statist.* **38** 191–200. [MR0207136](#)
- [10] Facer, M.R. and Müller, H.-G. (2003). Nonparametric estimation of the location of a maximum in a response surface. *J. Multivariate Anal.* **87** 191–217. [MR2007268](#)
- [11] Freedman, D. (1999). On the Bernstein–von Mises theorem with infinite-dimensional parameters. *Ann. Statist.* **27** 1119–1140. [MR1740119](#)
- [12] Ghosal, S. and van der Vaart, A.W. (2017). *Fundamentals of Nonparametric Bayesian Inference*. Cambridge: Cambridge Univ. Press. [MR3587782](#)

- [13] Hasminskii, R.Z. (1979). Lower bound for the risks of nonparametric estimates of the mode. In *Contributions to Statistics* (J. Jureckova, ed.) 91–97. Dordrecht: Reidel. [MR0561262](#)
- [14] Jørgensen, M., Nielsen, C.T., Keiding, N. and Skakkeback, N.E. (1985). Parametrische und nicht-parametrische Modelle für Wachstumsdaten. In *Neuere Verfahren der nichtparametrischen Statistik* (G.C. Pflug, ed.). *Med. Informatik und Statistik* **60** 74–87. Berlin: Springer.
- [15] Kiefer, J. and Wolfowitz, J. (1952). Stochastic estimation of the maximum of a regression function. *Ann. Math. Statist.* **23** 462–466. [MR0050243](#)
- [16] Kim, J. and Pollard, D. (1990). Cube root asymptotics. *Ann. Statist.* **18** 191–219. [MR1041391](#)
- [17] Knapik, B.T., van der Vaart, A.W. and van Zanten, J.H. (2011). Bayesian inverse problems with Gaussian priors. *Ann. Statist.* **39** 2626–2657. [MR2906881](#)
- [18] Lan, Y., Banerjee, M. and Michailidis, G. (2009). Change-point estimation under adaptive sampling. *Ann. Statist.* **37** 1752–1791. [MR2533471](#)
- [19] Mokkadem, A. and Pelletier, M. (2007). A companion for the Kiefer–Wolfowitz–Blum stochastic approximation algorithm. *Ann. Statist.* **35** 1749–1772. [MR2351104](#)
- [20] Müller, H.-G. (1985). Kernel estimators of zeros and of location and size of extrema of regression functions. *Scand. J. Stat.* **12** 221–232. [MR0817940](#)
- [21] Müller, H.-G. (1989). Adaptive nonparametric peak estimation. *Ann. Statist.* **17** 1053–1069. [MR1015137](#)
- [22] Polyak, B.T. and Tsybakov, A.B. (1990). Optimal order of accuracy of search algorithms in stochastic optimization. *Probl. Inf. Transm.* **26** 126–133. [MR1074128](#)
- [23] Schumaker, L.L. (2007). *Spline Functions: Basic Theory*, 3rd ed. New York: Cambridge Univ. Press. [MR2348176](#)
- [24] Shen, W. and Ghosal, S. (2015). Adaptive Bayesian procedures using random series priors. *Scand. J. Statist.* **42** 1194–1213. [MR3426318](#)
- [25] Shoung, J.-M. and Zhang, C.-H. (2001). Least squares estimators of the mode of a unimodal regression function. *Ann. Statist.* **29** 648–665. [MR1865335](#)
- [26] Silverman, B.W. (1985). Some aspects of the spline smoothing approach to nonparametric regression curve fitting (with discussion). *J. Roy. Statist. Soc. Ser. B* **47** 1–52. [MR0805063](#)
- [27] Szabó, B., van der Vaart, A.W. and van Zanten, J.H. (2015). Frequentist coverage of adaptive non-parametric Bayesian credible sets. *Ann. Statist.* **43** 1391–1428. [MR3357861](#)
- [28] Tang, R., Banerjee, M. and Michailidis, G. (2011). A two-stage hybrid procedure for estimating an inverse regression function. *Ann. Statist.* **39** 956–989. [MR2816344](#)
- [29] Tsybakov, A.B. (1990). Recursive estimation of the mode of a multivariate distribution. *Probl. Inf. Transm.* **26** 31–37. [MR1051586](#)
- [30] van der Vaart, A.W. (1998). *Asymptotic Statistics*. Cambridge: Cambridge Univ. Press. [MR1652247](#)
- [31] van der Vaart, A.W. and van Zanten, J.H. (2008). Reproducing kernel Hilbert spaces of Gaussian priors. In *Pushing the Limits of Contemporary Statistics: Contributions in Honor of Jayanta K. Ghosh. Inst. Math. Stat. (IMS) Collect.* **3** 200–222. Beachwood, OH: IMS. [MR2459226](#)
- [32] van der Vaart, A.W. and Wellner, J.A. (1996). *Weak Convergence and Empirical Processes With Applications to Statistics*. New York: Springer. [MR1385671](#)
- [33] Yoo, W.W. and Ghosal, S. (2016). Supremum norm posterior contraction and credible sets for non-parametric multivariate regression. *Ann. Statist.* **44** 1069–1102. [MR3485954](#)
- [34] Yoo, W.W. and Ghosal, S. (2018). Supplement to “Bayesian mode and maximum estimation and accelerated rates of contraction.” DOI:10.3150/18-BEJ1056SUPP.
- [35] Yoo, W.W., Rivoirard, V. and Rousseau, J. (2018). Adaptive supremum norm posterior contraction: Wavelet spike-and-slab and anisotropic Besov spaces. [arXiv:1708.01909 \[math.ST\]](#).
- [36] Yoo, W.W. and van der Vaart, A.W. (2018). The Bayes Lepski’s method and credible bands through volume of tubular neighborhoods. [arXiv:1711.06926 \[math.ST\]](#).

Received August 2017 and revised March 2018